

# Повышаем живучесть Raft в реальных условиях



Сергей  
Останевич



**HighLoad<sup>++</sup>**  
2022



Tarantool

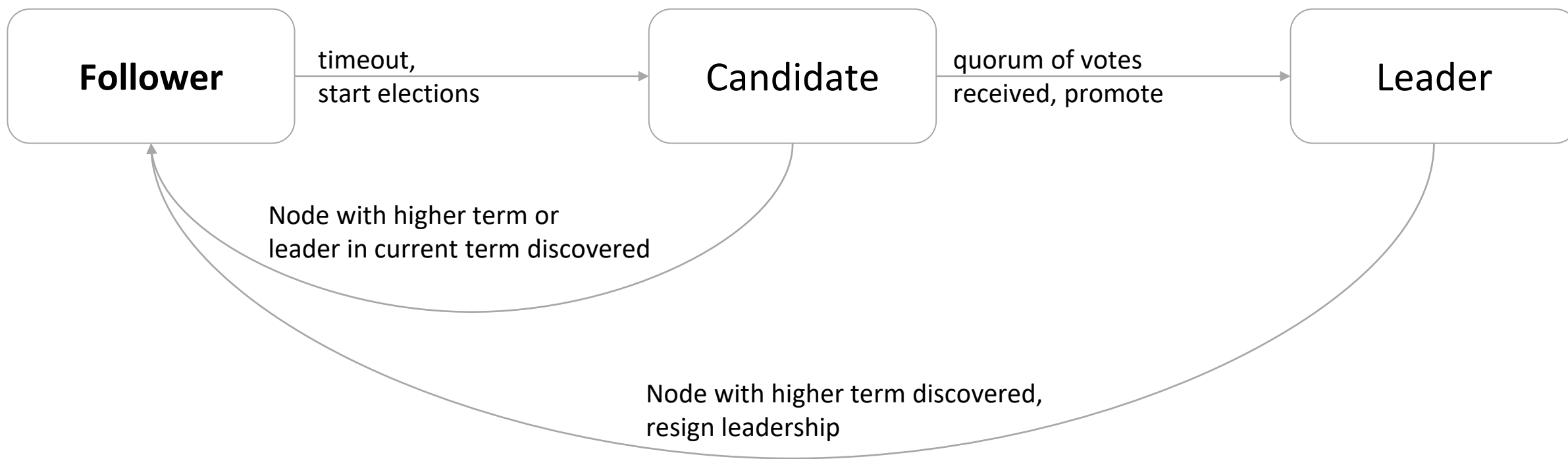
# Меню

- **Raft overview**
  - Термины: Journal, Term, LSN
  - Выборы
  - Гарантии
  - Ожидания != Реальность
- Raft / Tarantool: особенности
- Настройки Raft
  - Pre-Vote
  - Split-Vote detection
  - Fencing

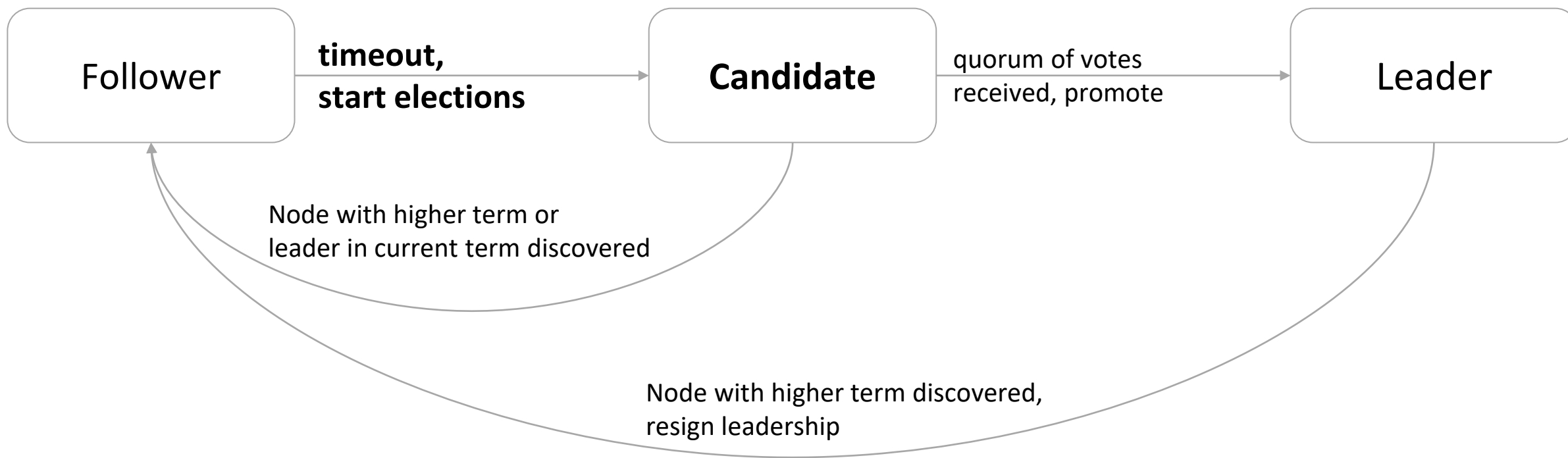
# Raft overview: Термины

- **Raft** – алгоритм достижения консенсуса в распределенных системах
- Raft достигает консенсуса за счет **выбора единого лидера**, который может изменять состояние системы
- Лидер выбирается в рамках **term** – пронумерованного отрезка времени
- Лидер может добавлять записи в **журнал (log)**
- Каждая запись в журнале пронумерована **log index/log sequence number (LSN)**

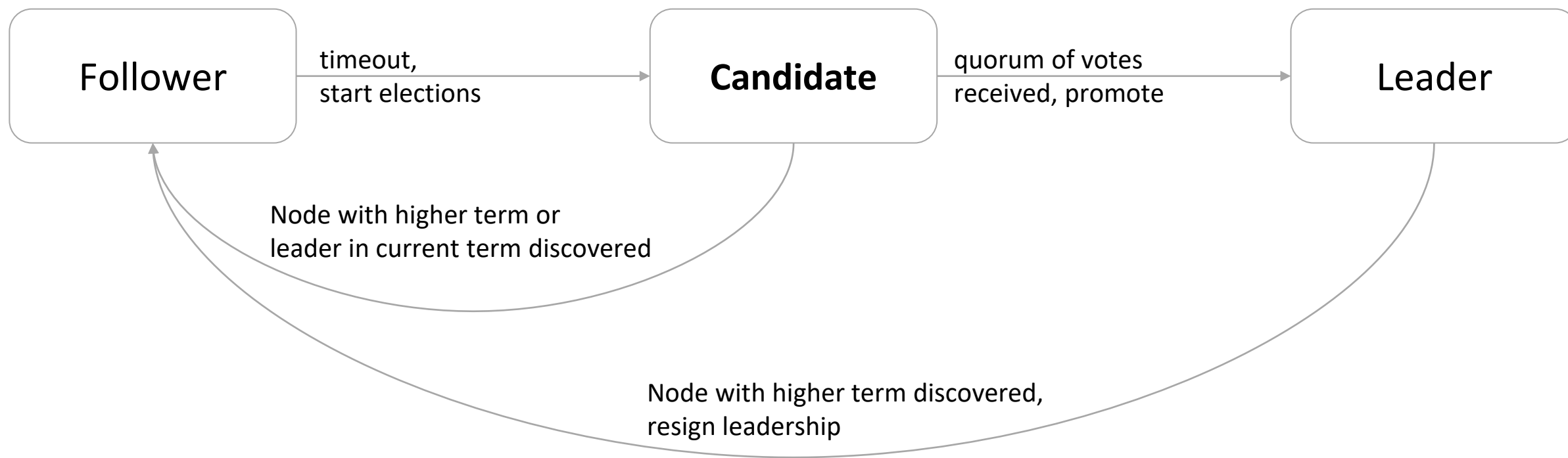
# Raft overview: Выборы



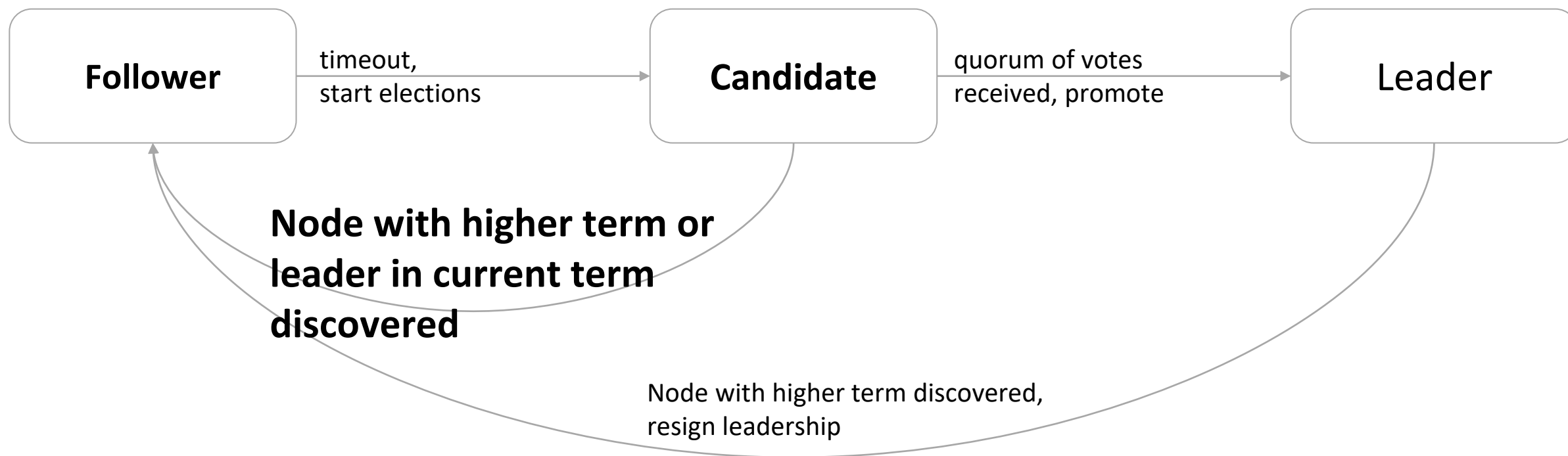
# Raft overview: Выборы



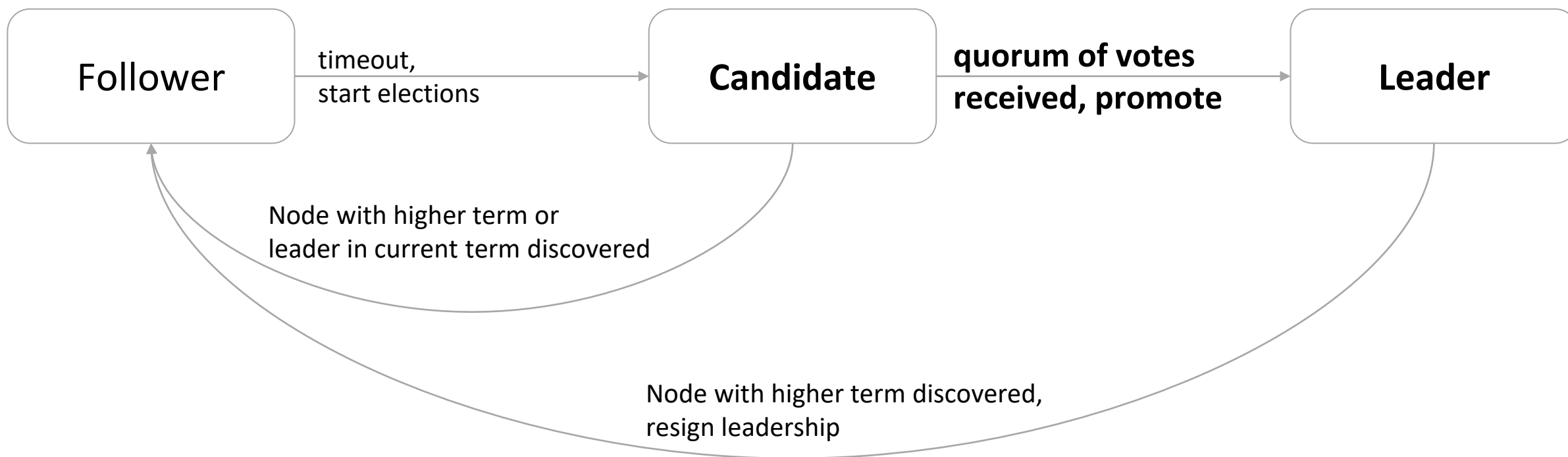
# Raft overview: Выборы



# Raft overview: Выборы

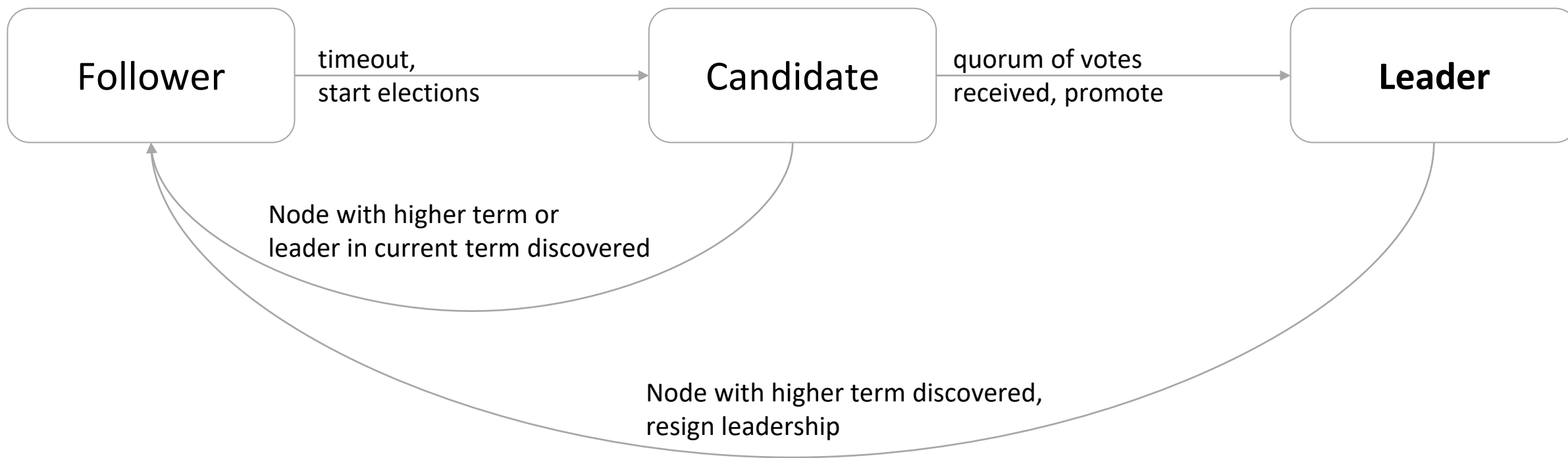


# Raft overview: Выборы

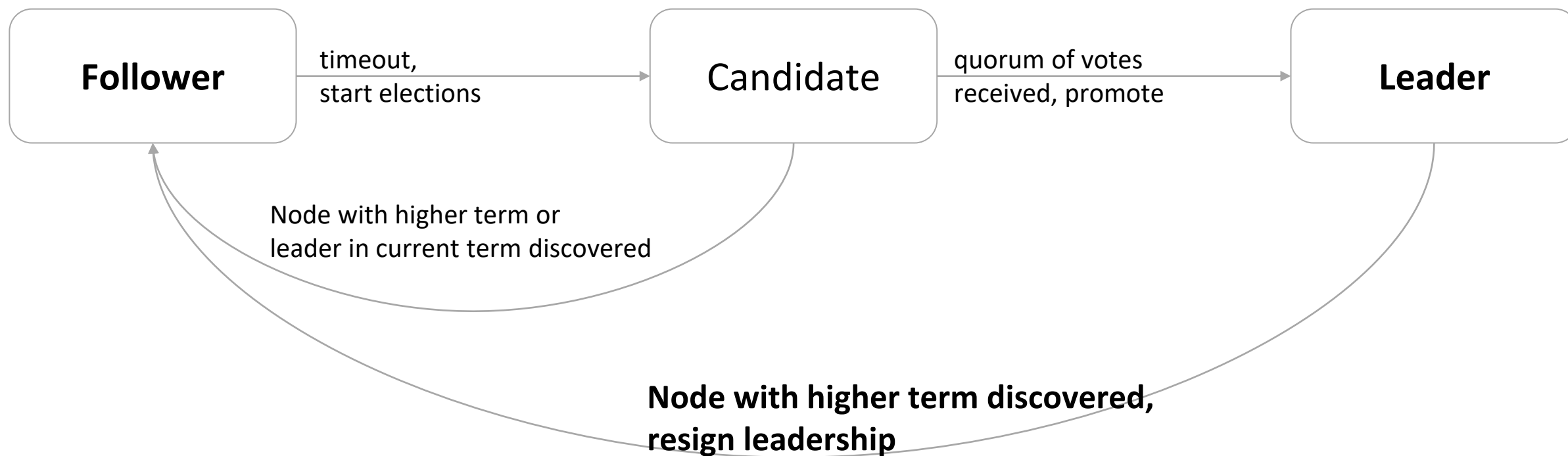




# Raft overview: Выборы



# Raft overview: Выборы



# Raft overview: Гарантии

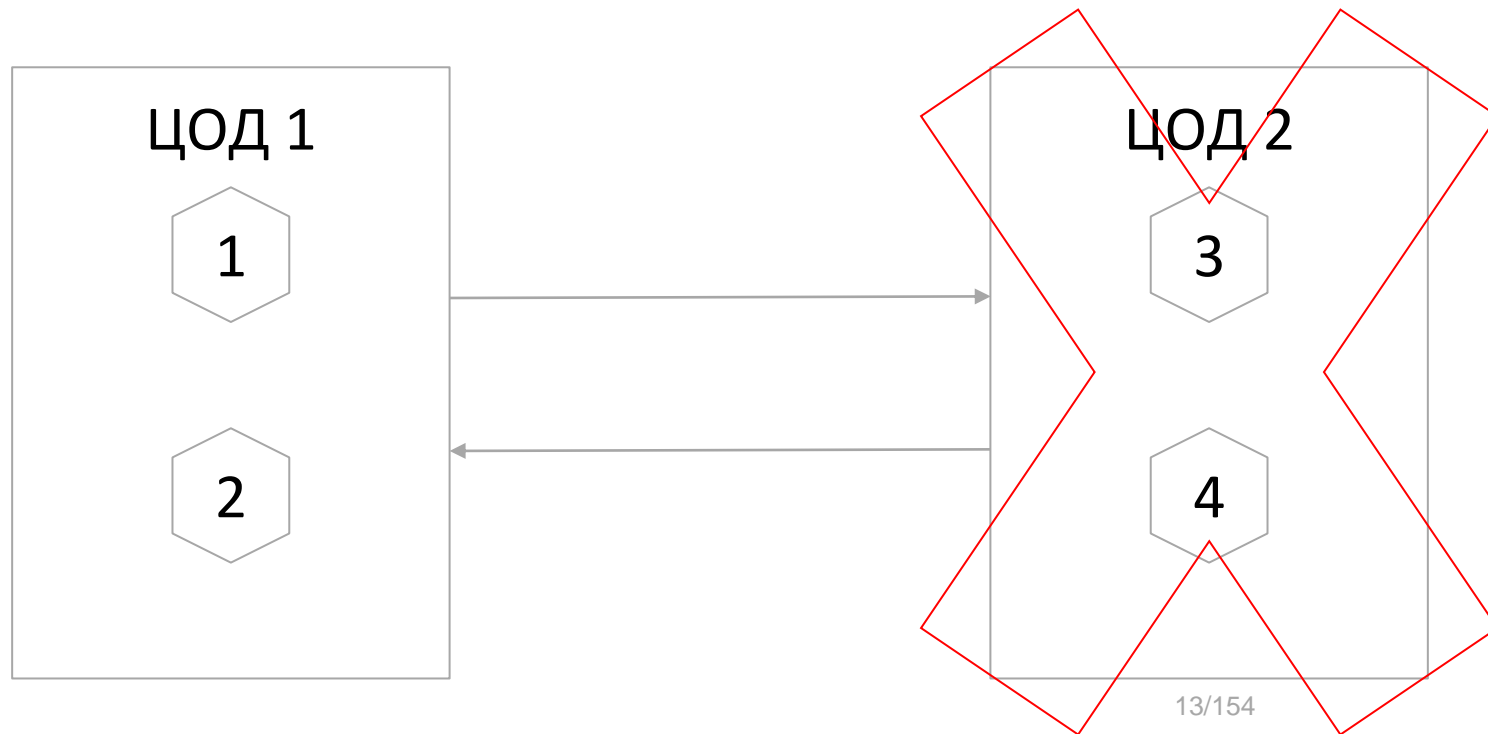
- **Election Safety:** один лидер в каждом term
- **Leader Append-Only:** лидер только добавляет записи без модификации или удаления записей
- **Log Matching:** записи с идентичными индексом и term в двух журналах гарантируют их совпадение вплоть до этих записей
- **Leader Completeness:** запись в журнал в каком-либо term будет присутствовать в логах всех последующих лидеров
- **State Machine Safety:** запись, примененная с определенным индексом на одном сервере, будет такой же и на другом сервере

# Raft overview: Гарантии

- **Рано или поздно** лидер найдётся (скорее всего)
- **В рамках одного term** есть не более одного лидера

# Raft overview: Гарантии

- **Рано или поздно** лидер найдётся (скорее всего)
- **В рамках одного term** есть не более одного лидера



# Raft overview: Гарантии

- **Рано или поздно** лидер найдётся (скорее всего)
- В рамках одного **term** есть не более одного лидера

## При ограничениях

- В кластере есть нода, имеющая необходимое количество соединений, чтобы выиграть выборы

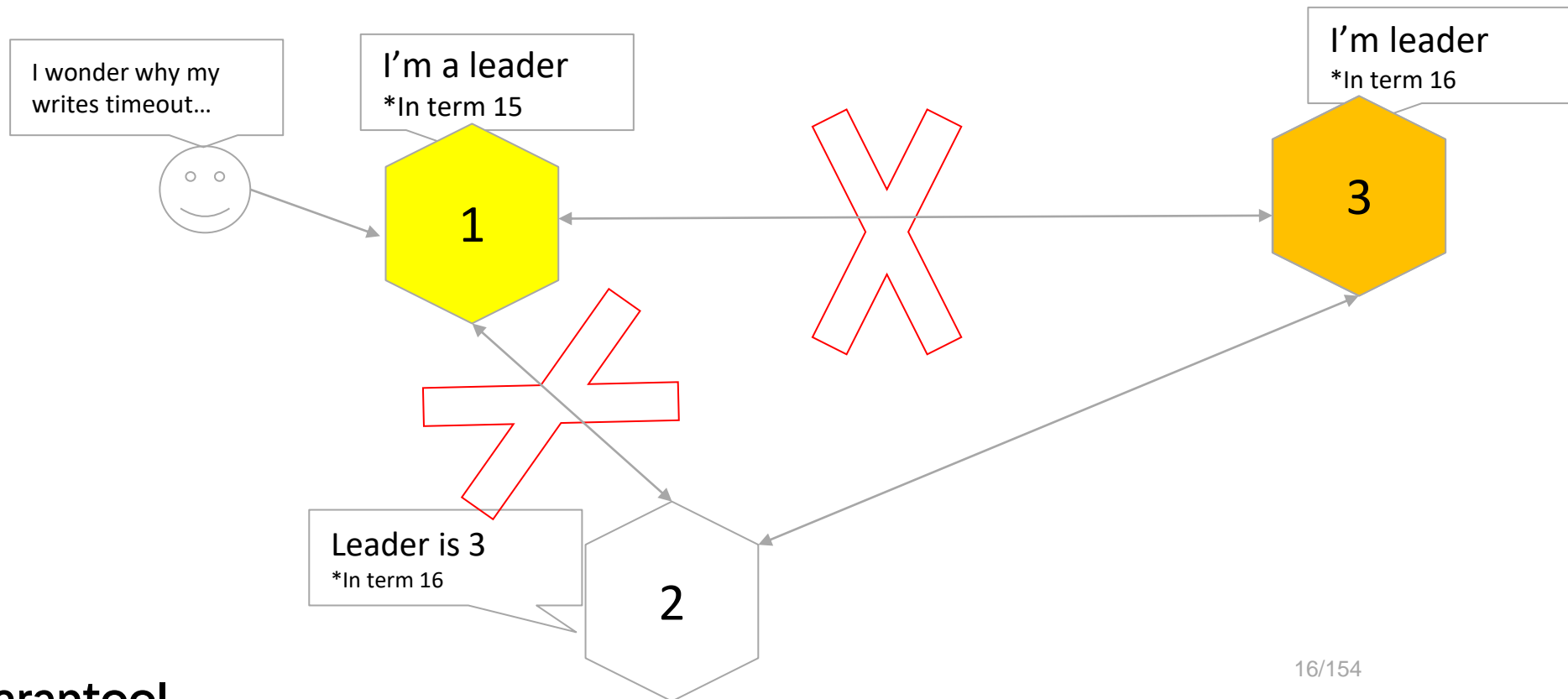
$$T_{ping} \ll T_{election} \ll T_{accident}$$

# Raft overview: Ожидания != Реальность

- Рано или поздно != быстро
- Term != время

# Raft overview: Ожидания != Реальность

- Рано или поздно != быстро
- Term != время





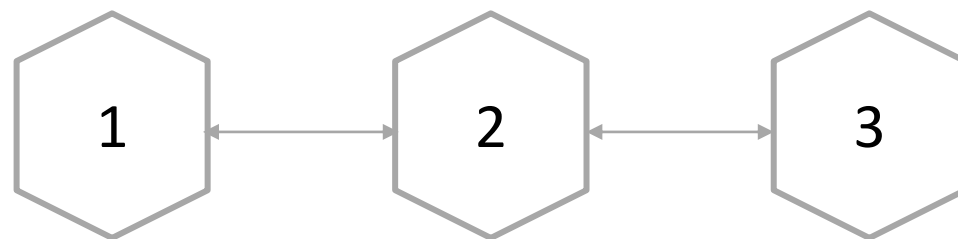
# Raft overview: Ожидания != Реальность

Ожидаем:

- Найти лидера **быстро**
- Не более одного лидера в кластере в **конкретный момент времени**
- Возвращение/подключение нод в кластер не вызывает его простоя

# Raft overview: Ожидания != Реальность

В процессе тестирования Raft в Tarantool мы нашли проблему в такой конфигурации:



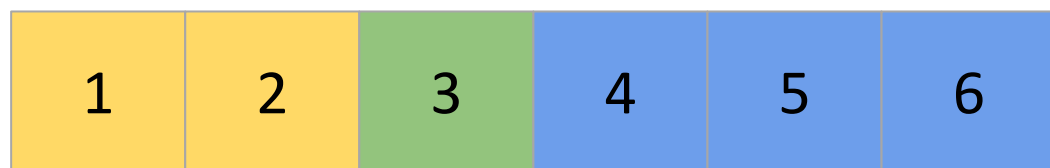
И почти сразу нашли, что мы не одиноки

The screenshot shows the header of a Cloudflare blog post. At the top left is the Cloudflare logo and the text 'The Cloudflare Blog'. To the right is a subscription form with the text 'Subscribe to receive notifications of new posts:', an 'Email Address' input field, and an orange 'Subscribe' button. Below the header is a navigation menu with links: 'Product News', 'Speed & Reliability', 'Security', 'Serverless', 'Zero Trust', 'Developers', 'Deep Dive', and 'Life @Cloudflare'. The main title of the post is 'A Byzantine failure in the real world' in a large, bold, dark font. Below the title is the date '27.11.2020'. The first line of the post body reads: 'On November 2, 2020, Cloudflare had an [incident](#) that impacted the availability of the API and dashboard for six hours and 33 minutes. During this incident, the'.

# Вторая переменная блюд

- ✓ Raft overview
  - Journal, Term, LSN
  - Выборы
  - Гарантии
  - Ожидания != Реальность
- **Raft / Tarantool: особенности**
- Настройки Raft
  - Pre-Vote
  - Split-Vote detection
  - Fencing

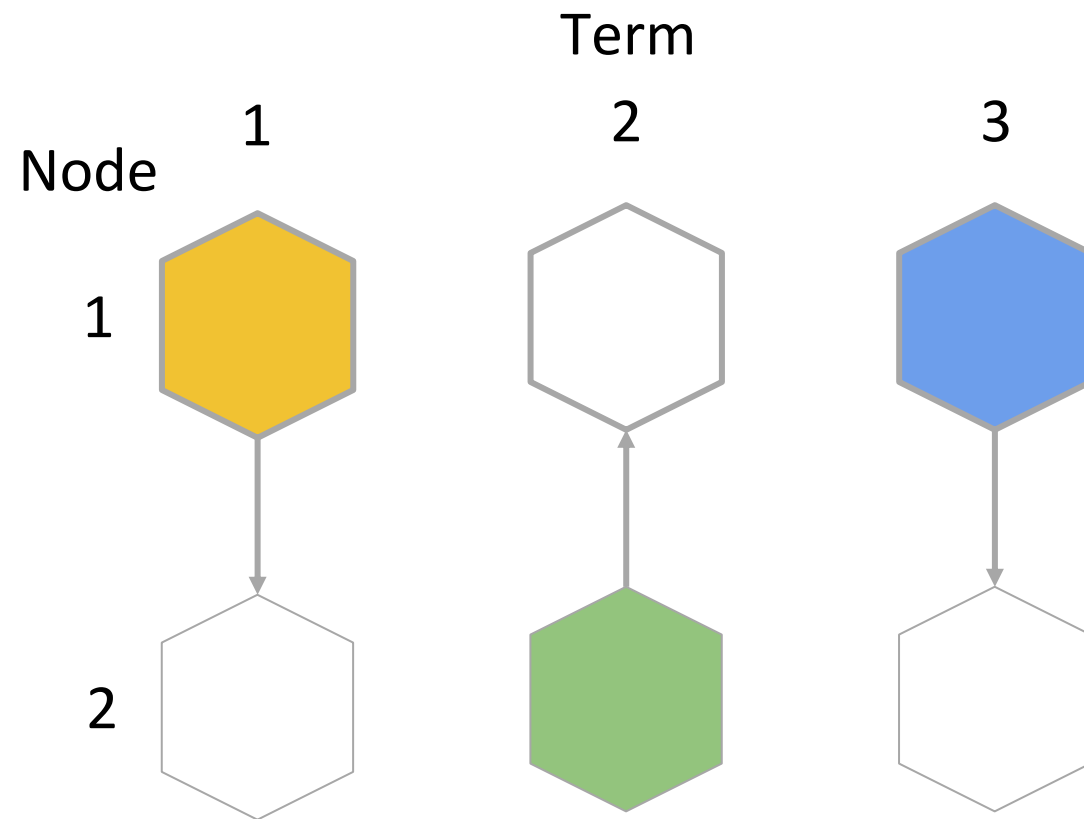
# Tarantool: Мультиплексированный журнал



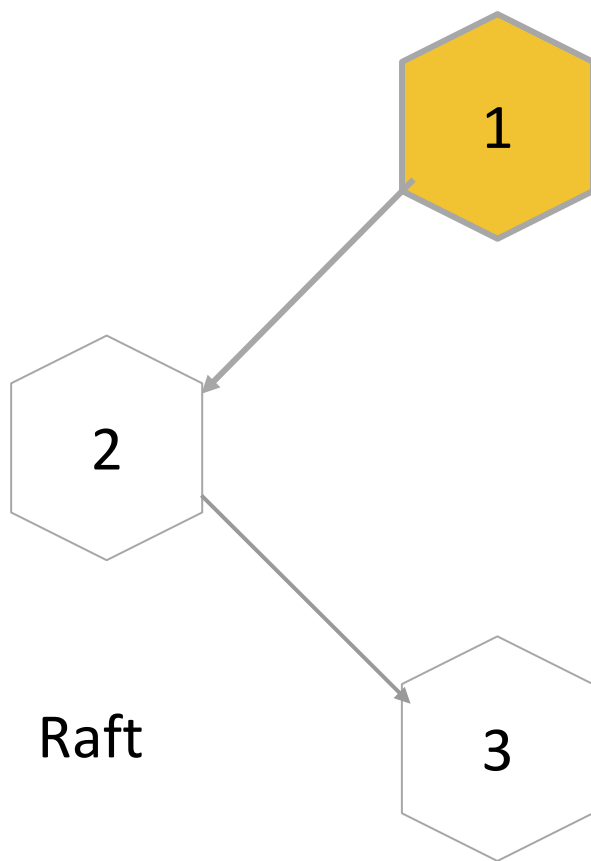
log index = 6



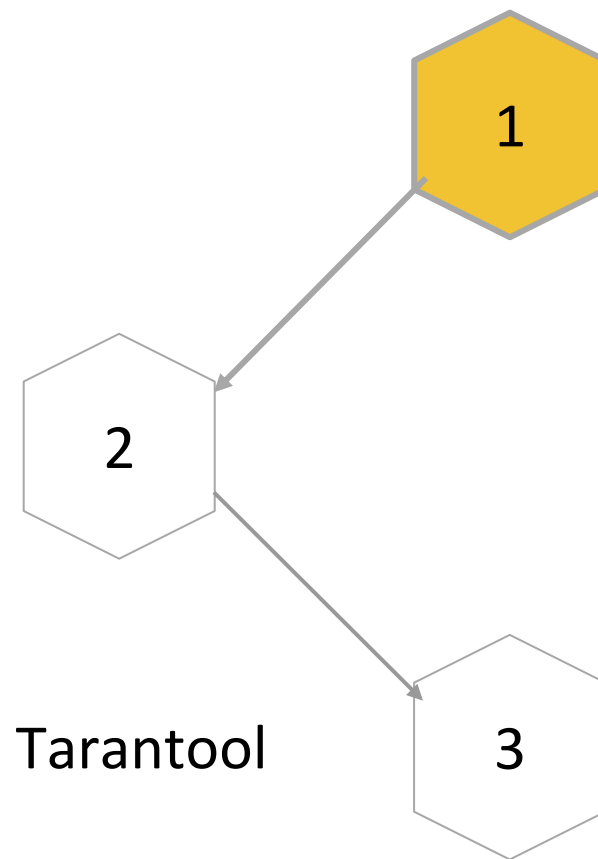
vclock = {1 : 5, 2: 1}



# Tarantool: Рассылка данных от всех ко всем

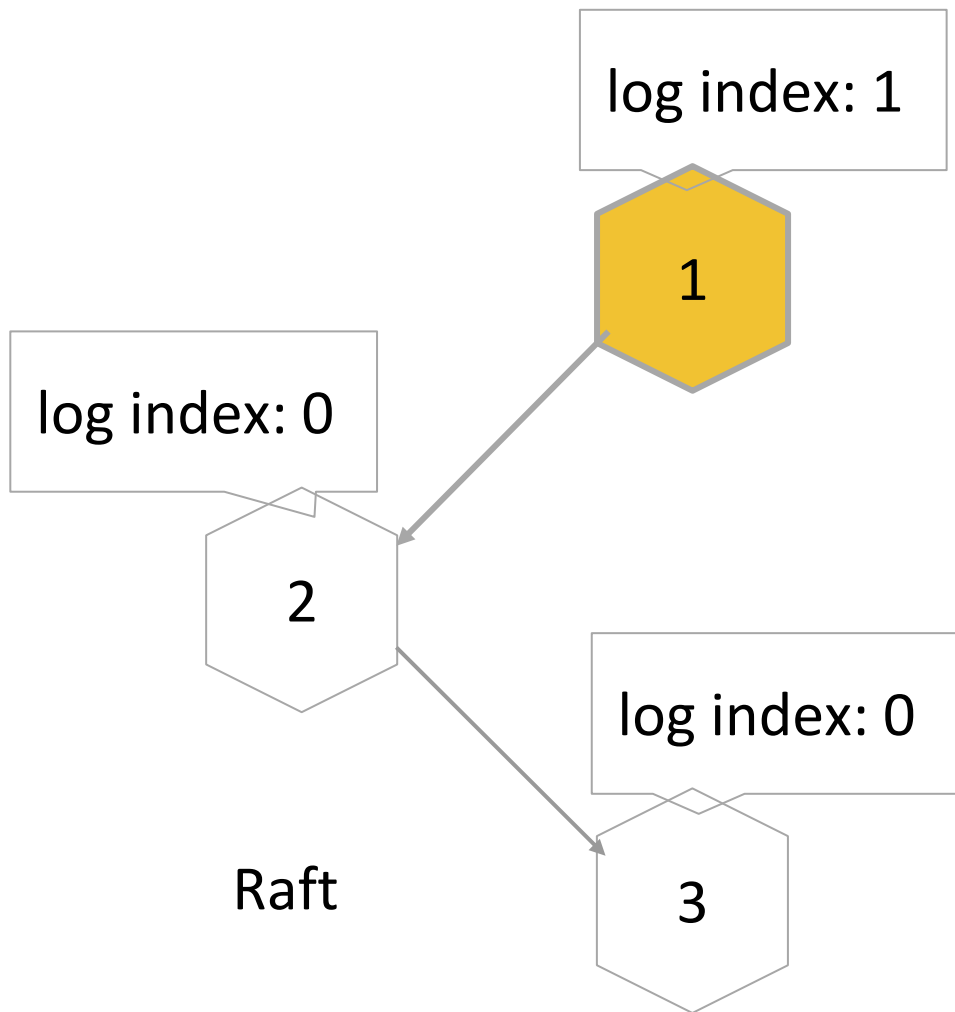


Raft

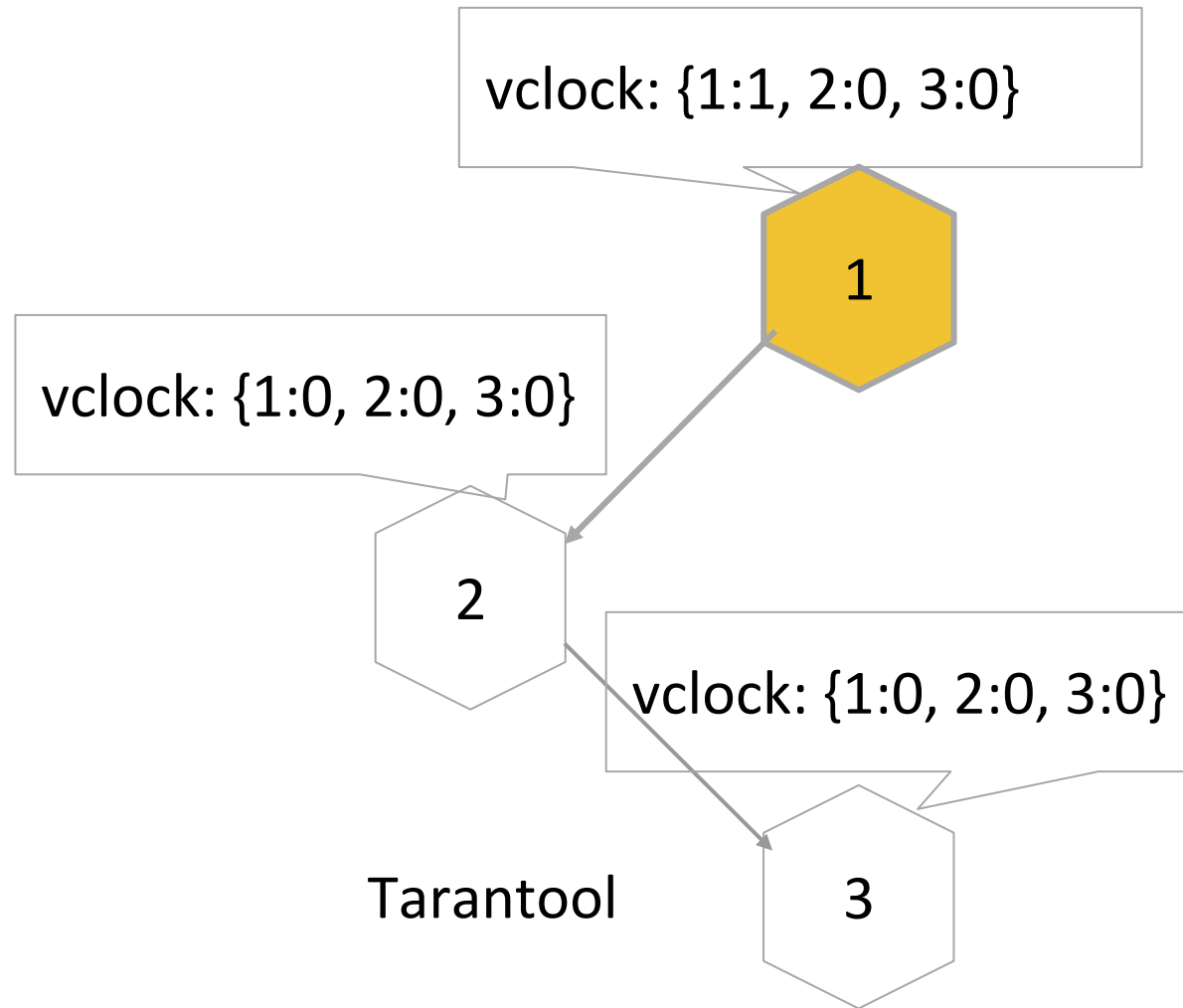


Tarantool

# Tarantool: Рассылка данных от всех ко всем

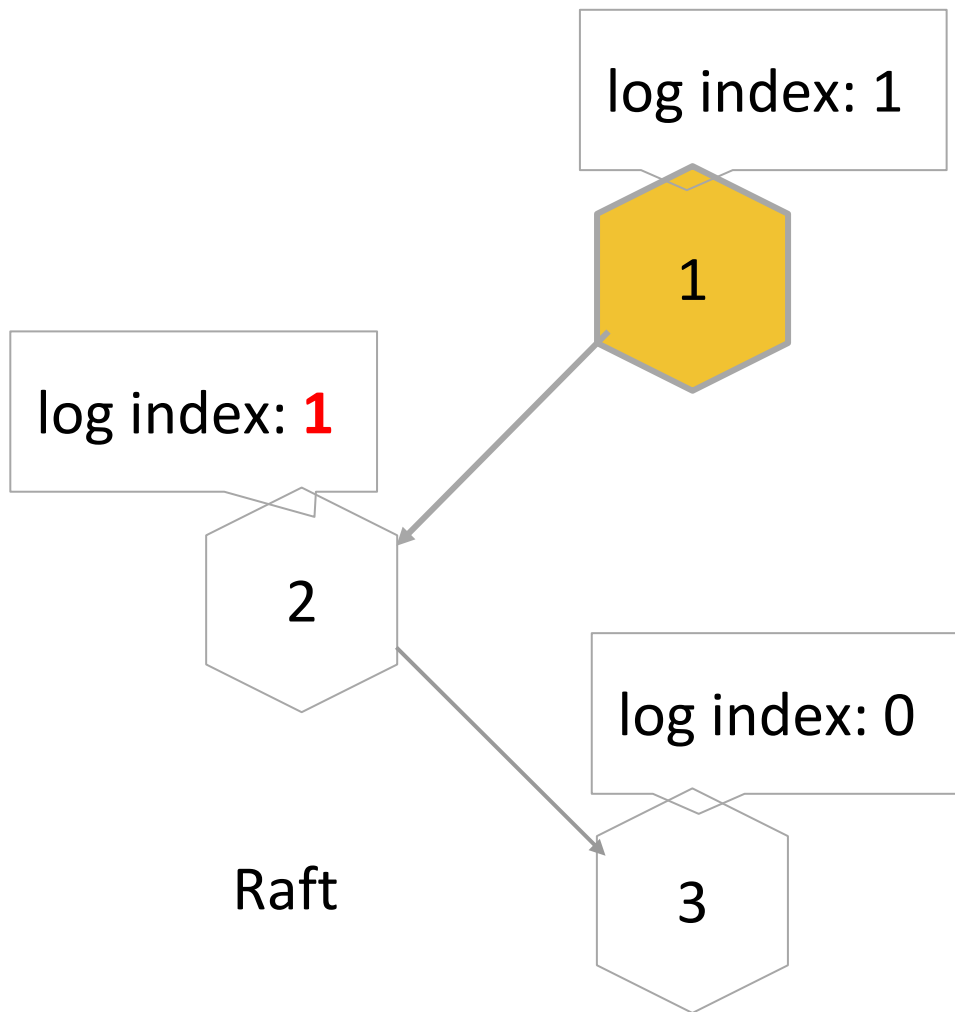


Raft

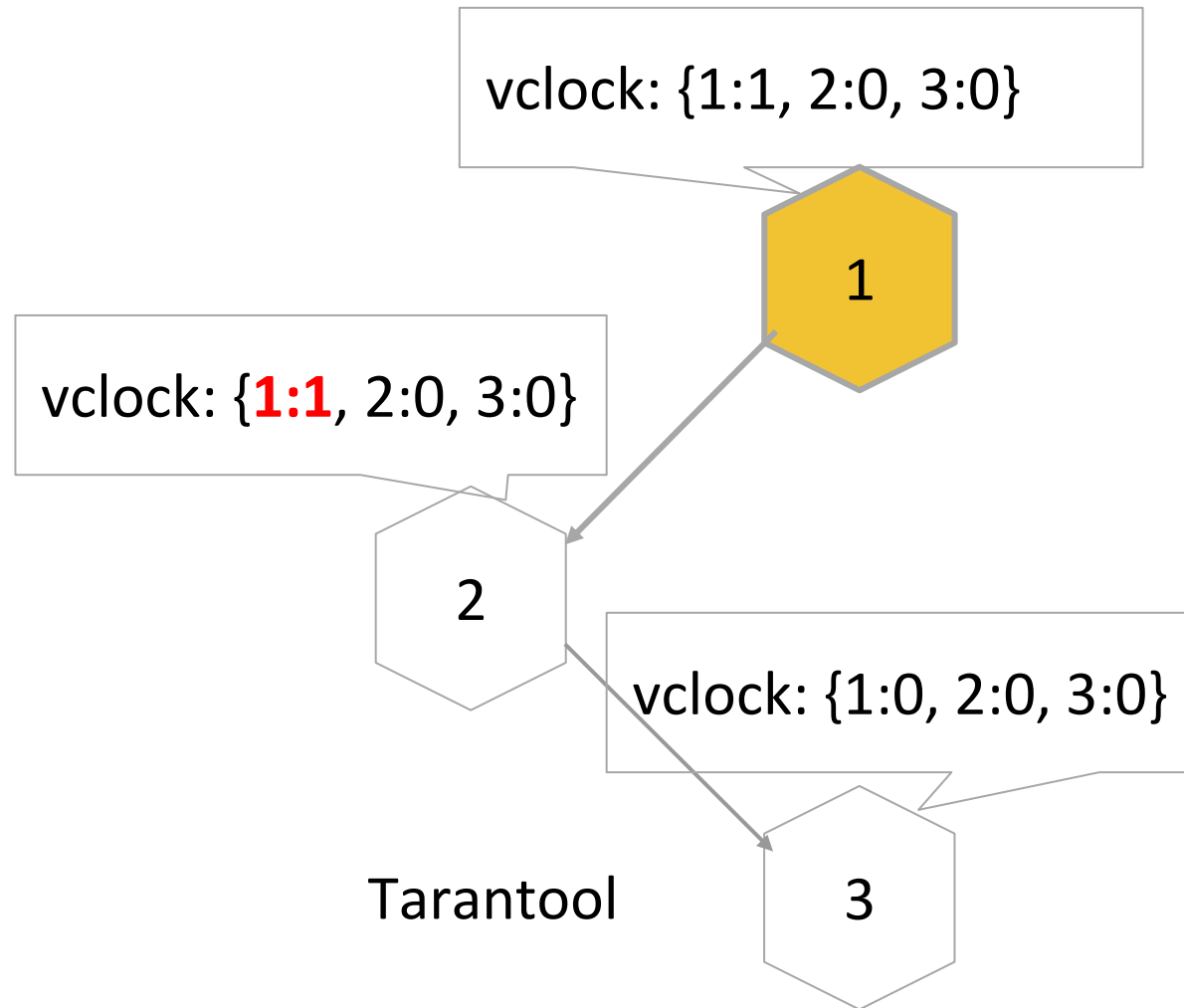


Tarantool

# Tarantool: Рассылка данных от всех ко всем

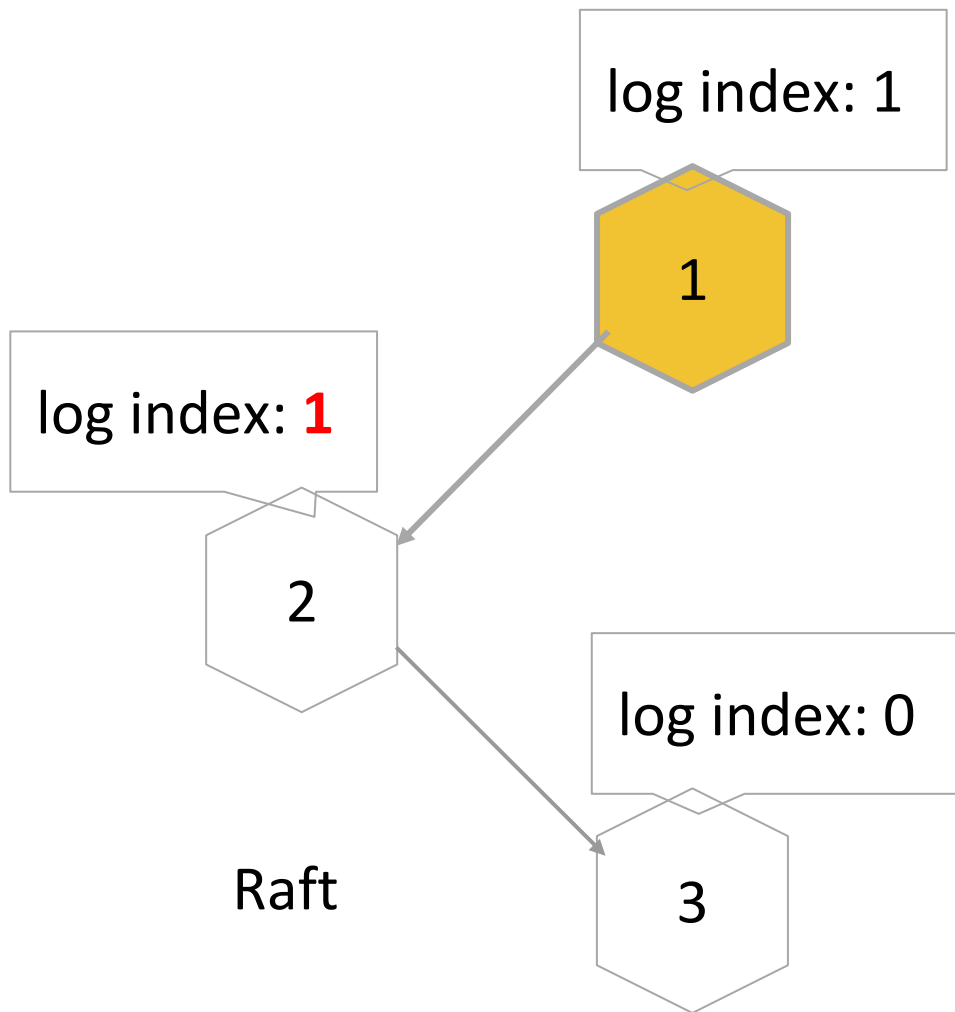


Raft

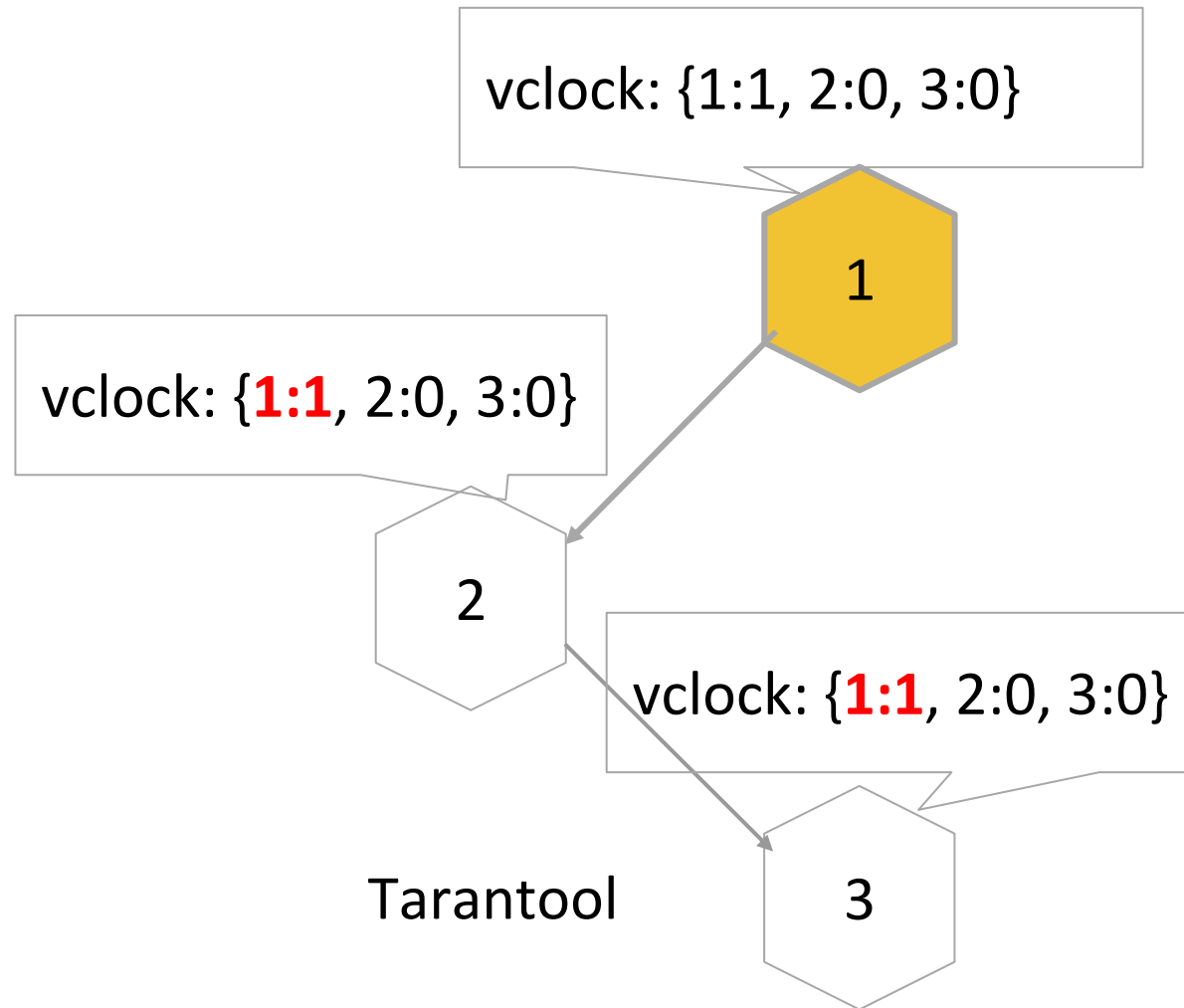


Tarantool

# Tarantool: Рассылка данных от всех ко всем



Raft

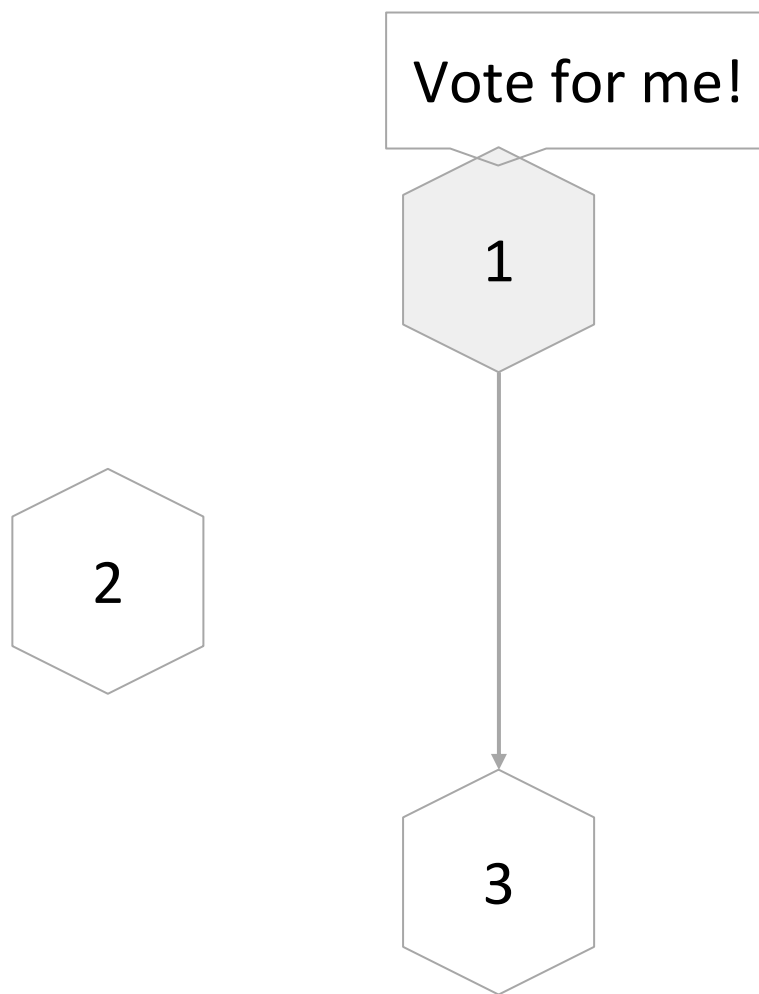


Tarantool

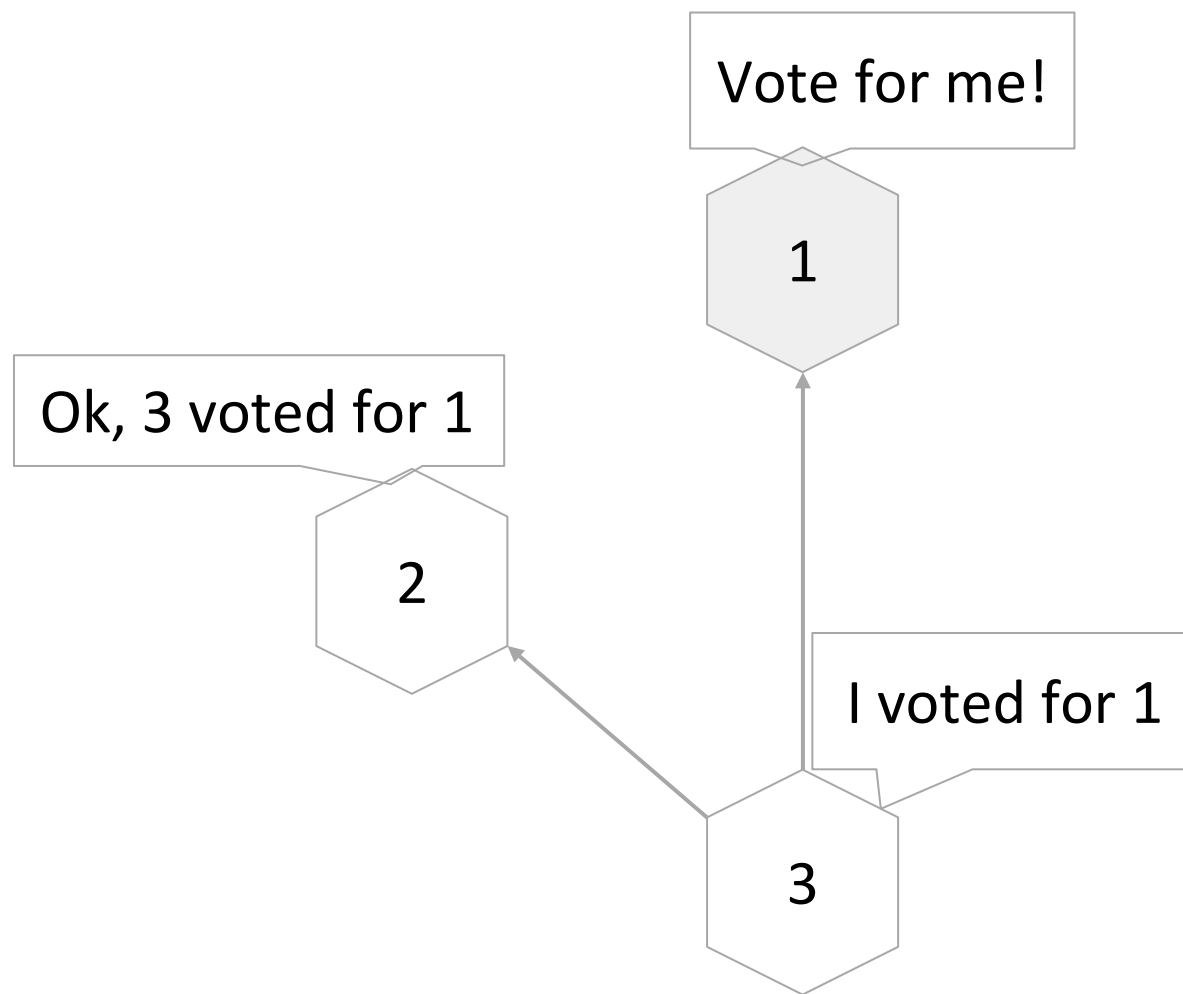




# Tarantool: БROADCAST результатов голосования



# Tarantool: Бroadcast результатов голосования



Сообщения существующих типов можно дополнять без нарушения обратной совместимости

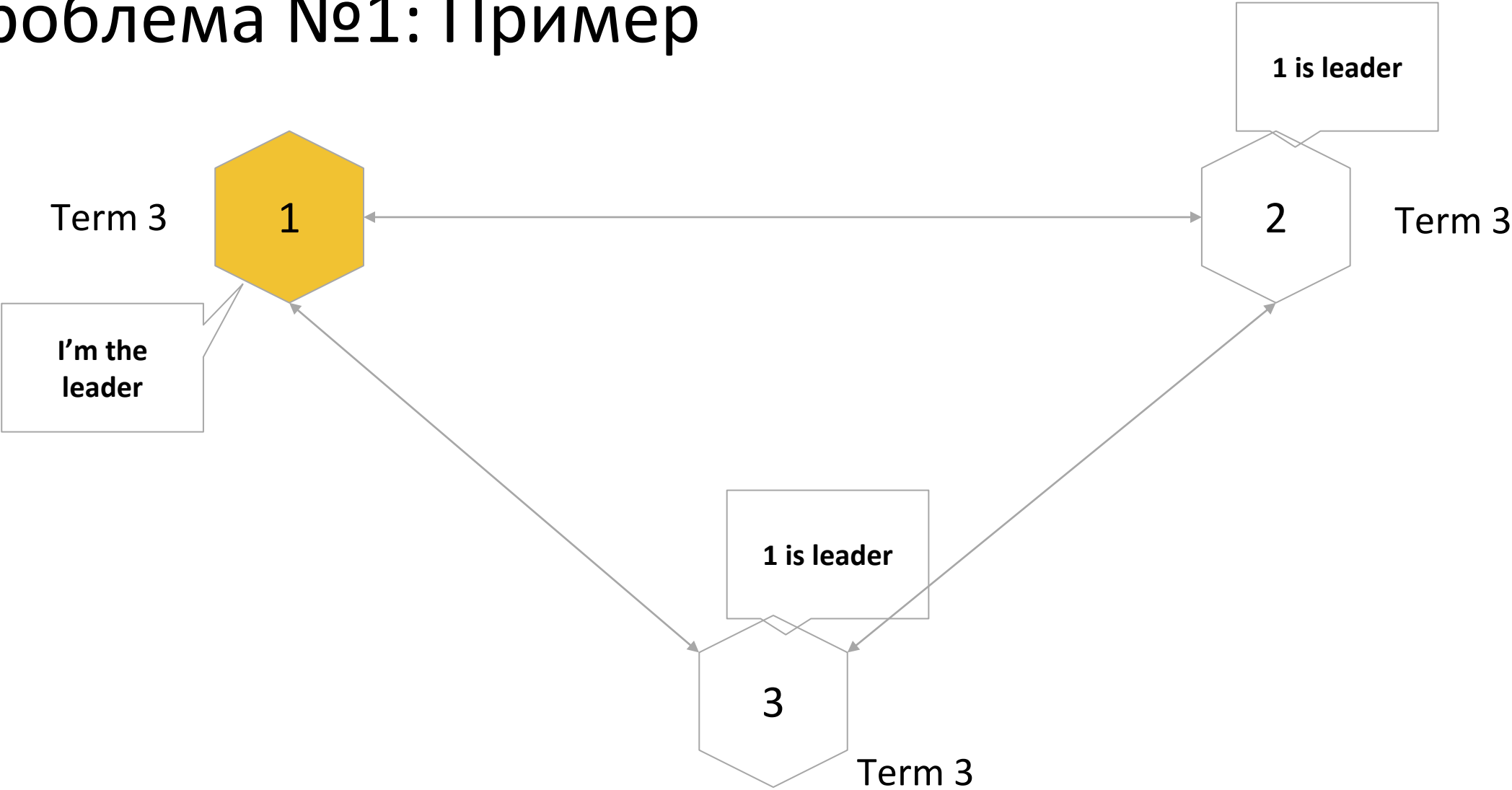
# Меню

- ✓ Raft overview
  - Термины: Journal, Term, LSN
  - Выборы
  - Гарантии
  - Ожидания != Реальность
- ✓ Raft / Tarantool: особенности
  - **Надстройки Raft**
    - **Pre-Vote**
    - Split-Vote detection
    - Fencing

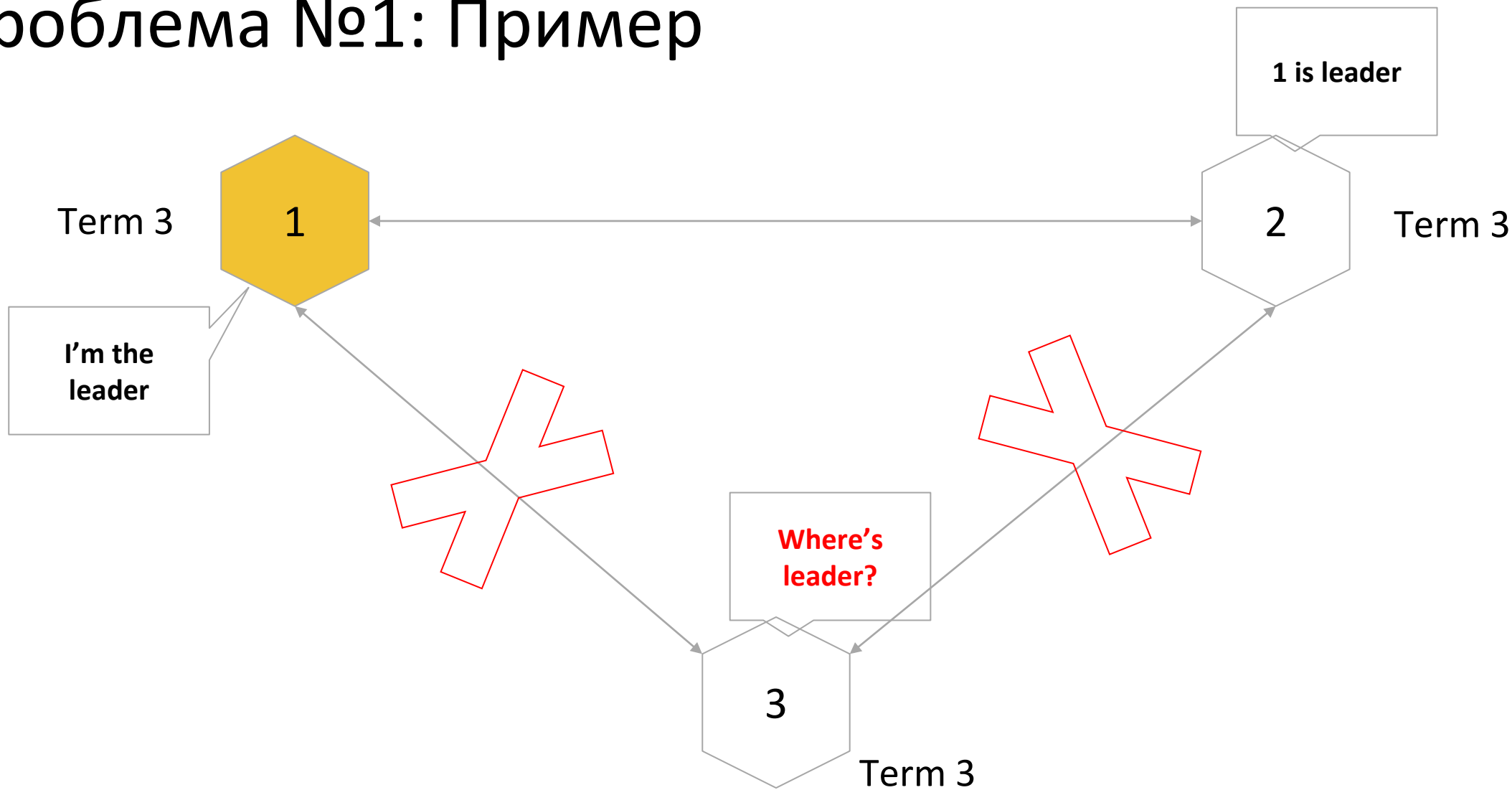
# Pre-Vote: Проблема №1

Недоступность кластера на запись после восстановления связности

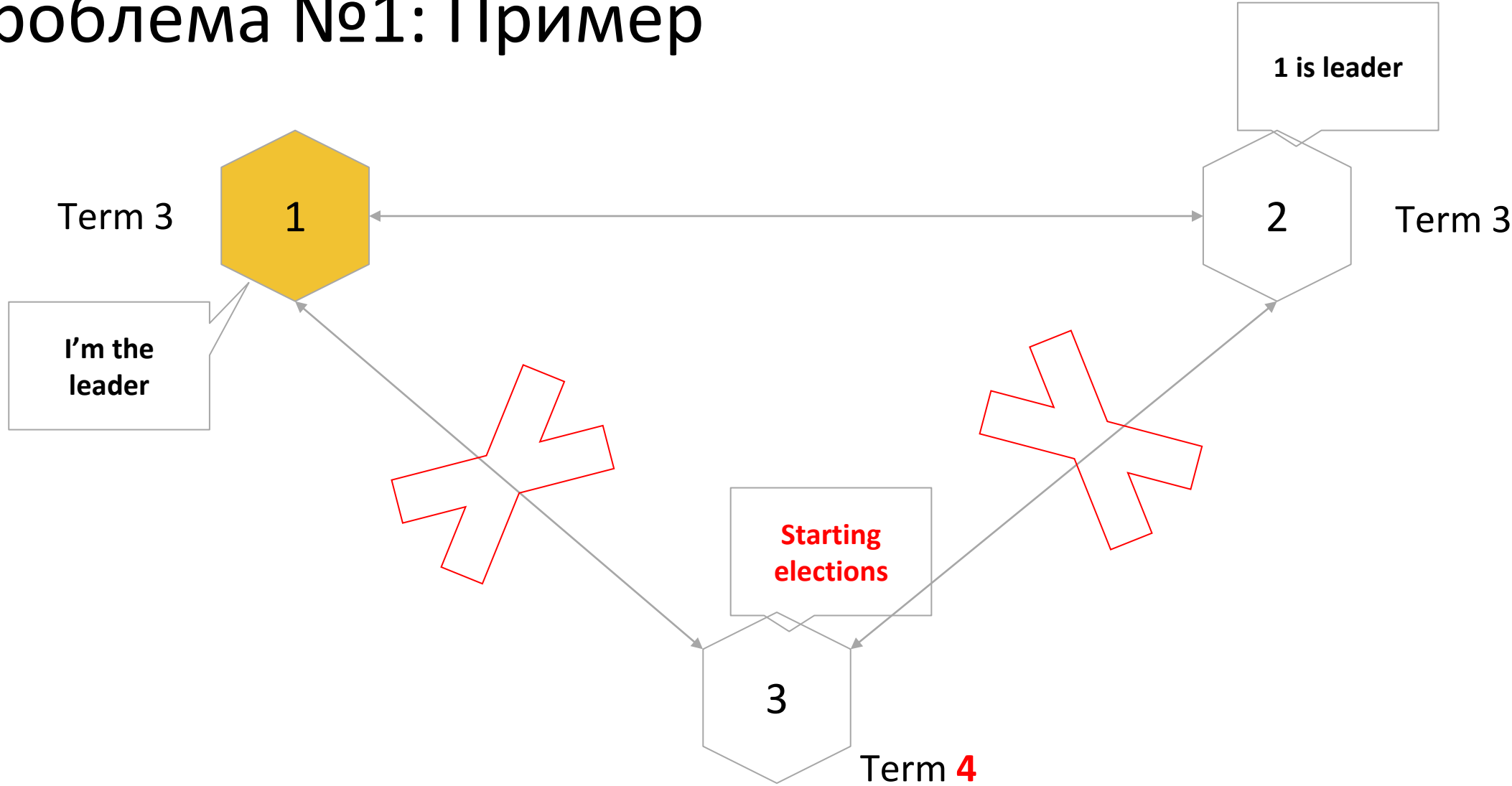
# Проблема №1: Пример



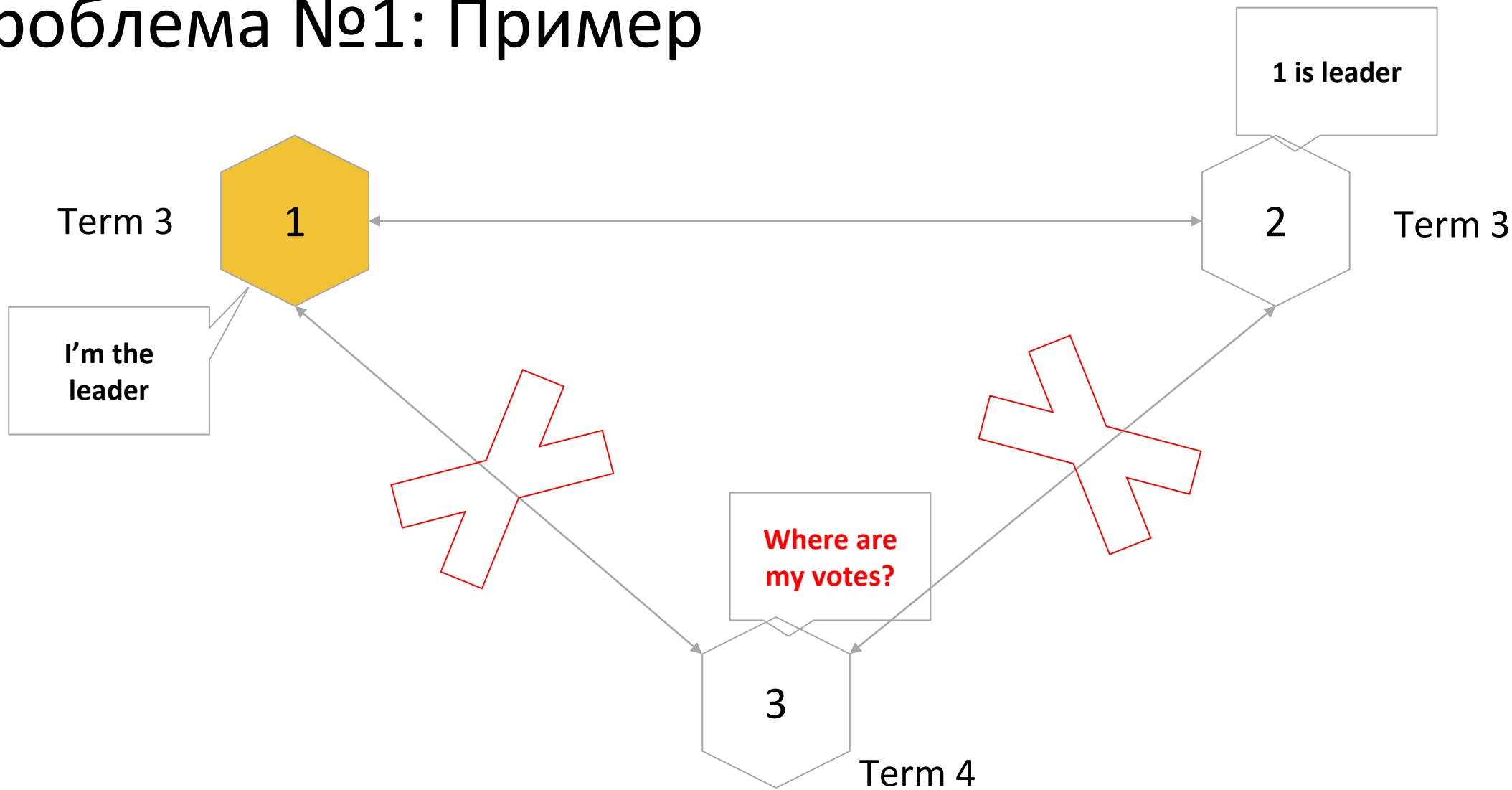
# Проблема №1: Пример



# Проблема №1: Пример

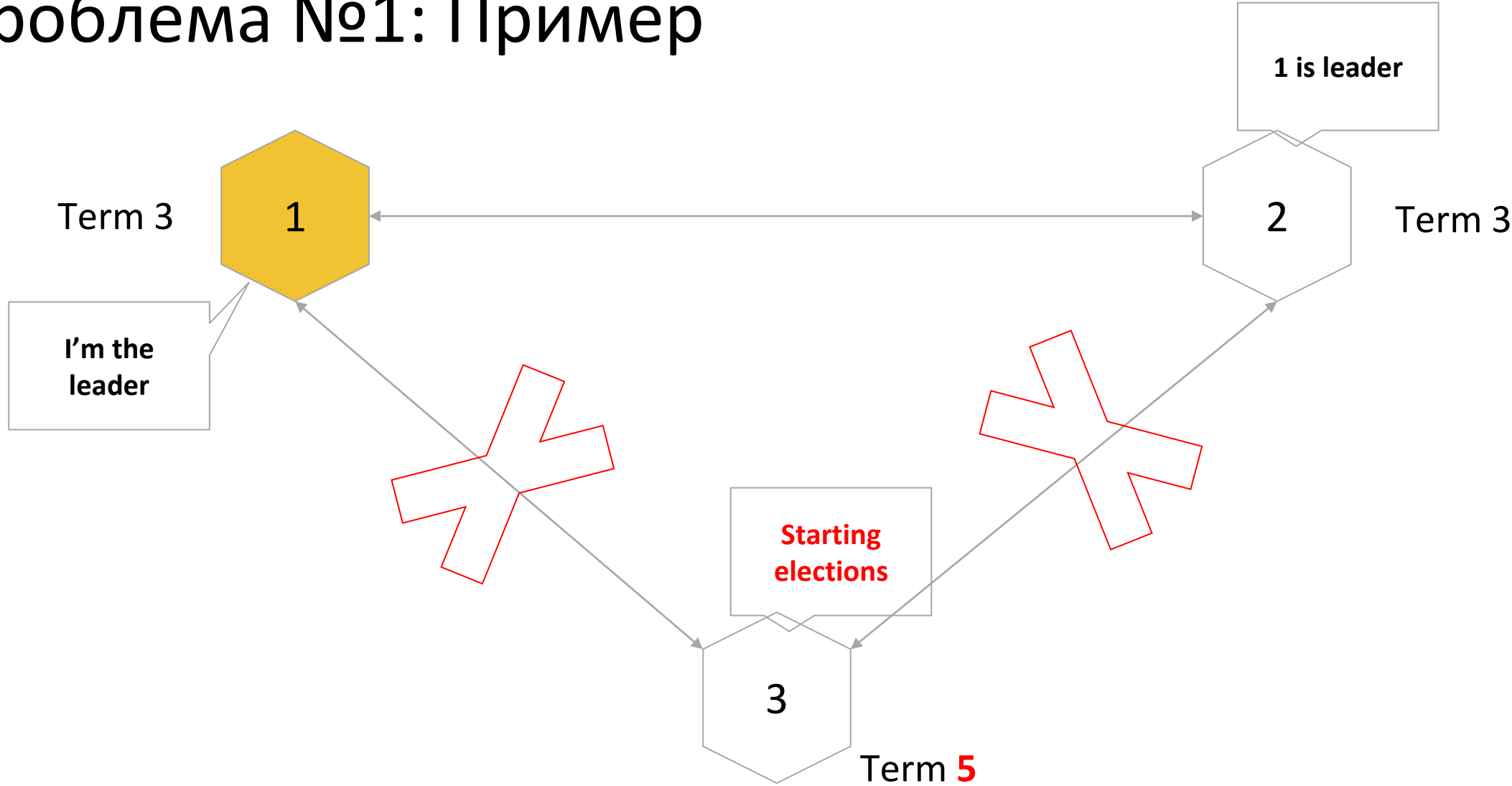


# Проблема №1: Пример

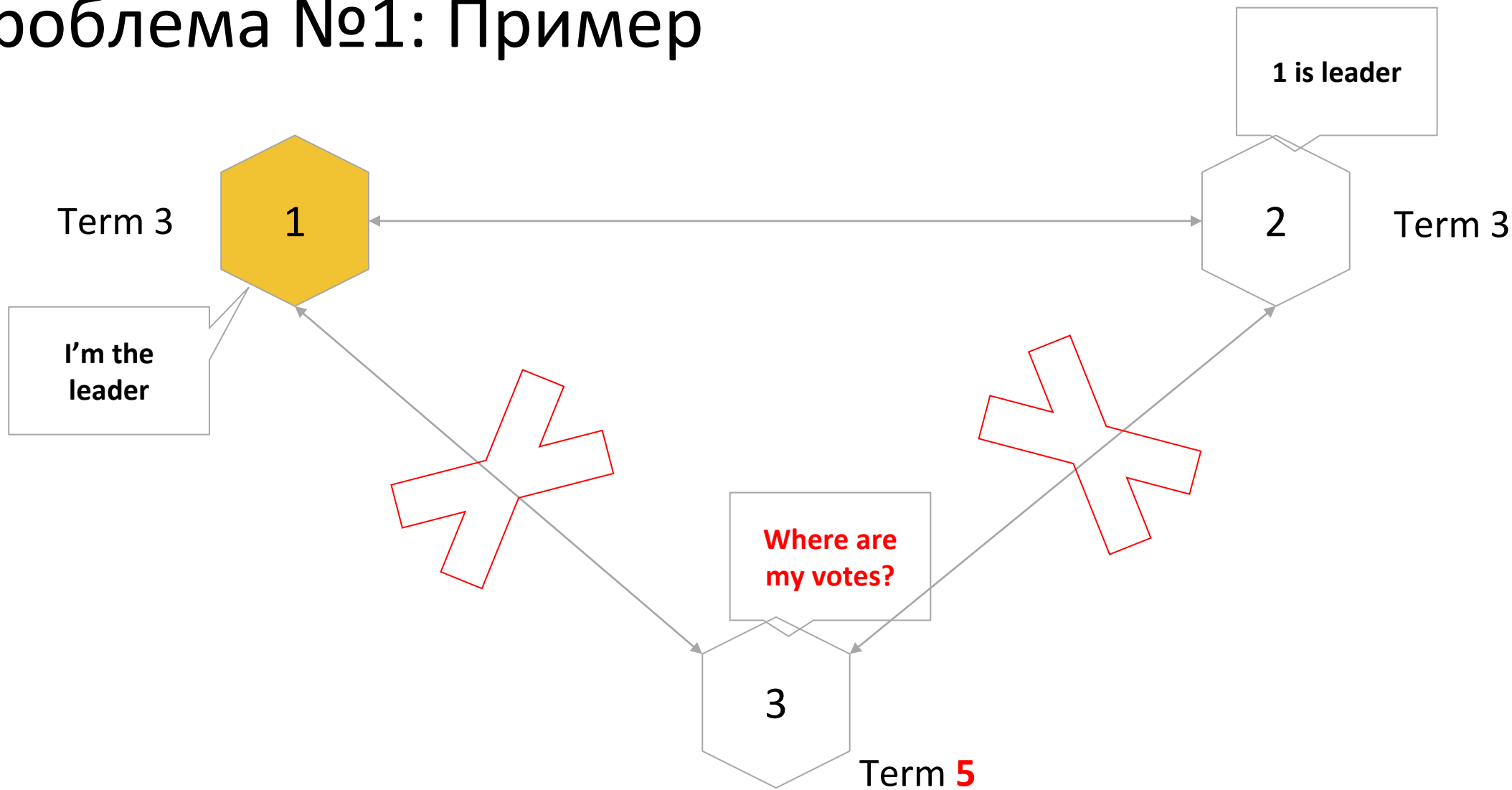




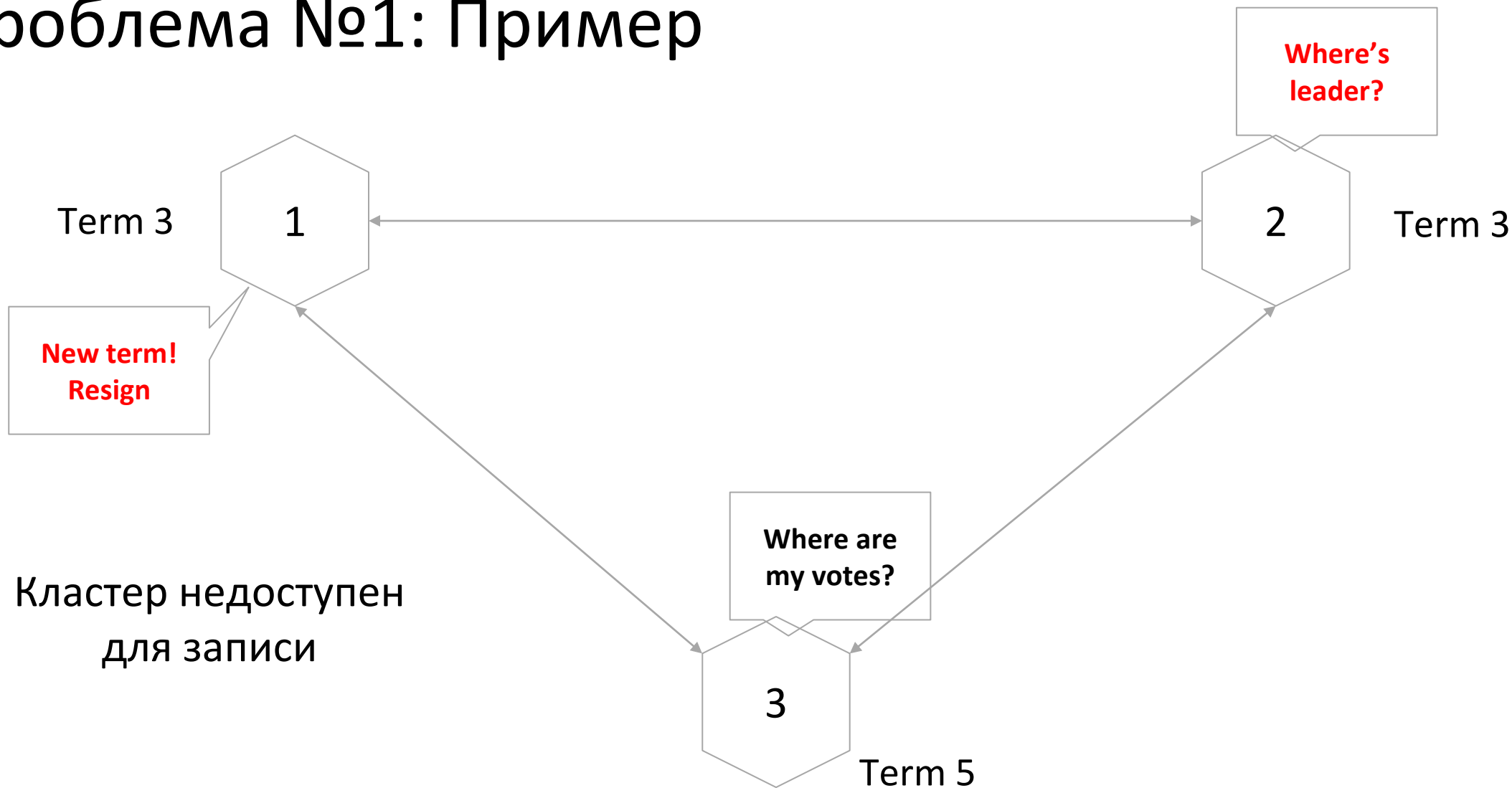
# Проблема №1: Пример



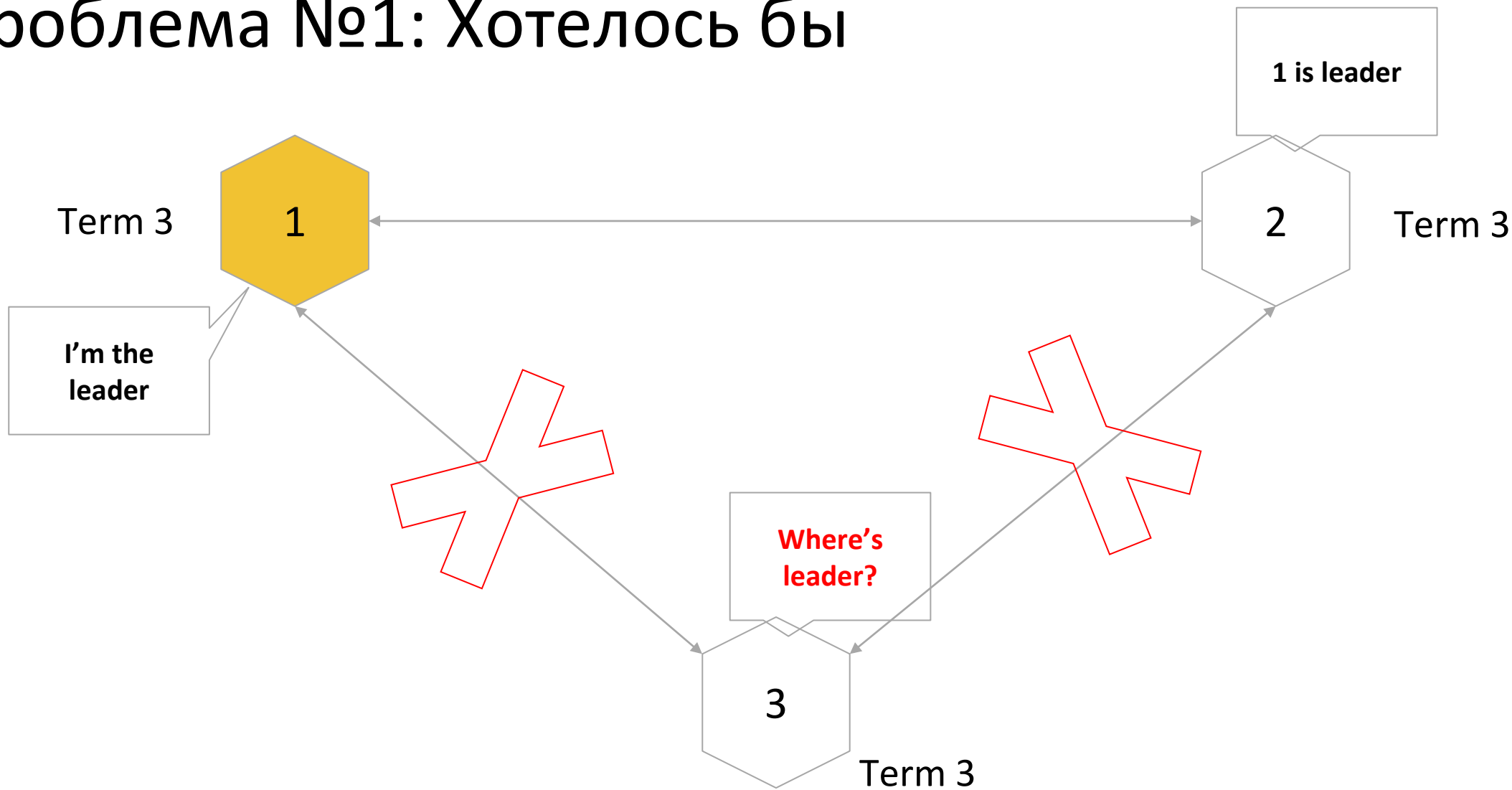
# Проблема №1: Пример



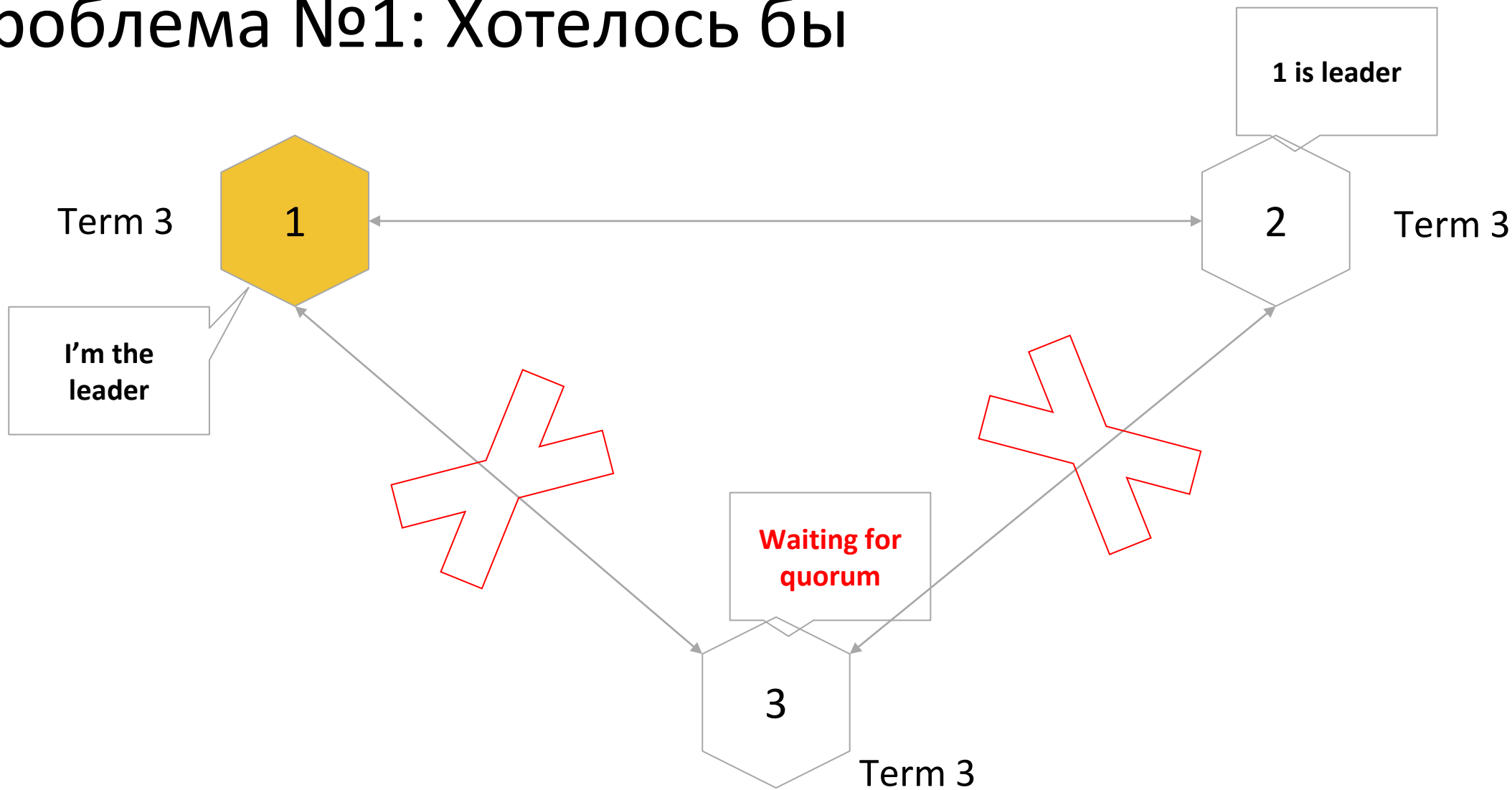
# Проблема №1: Пример



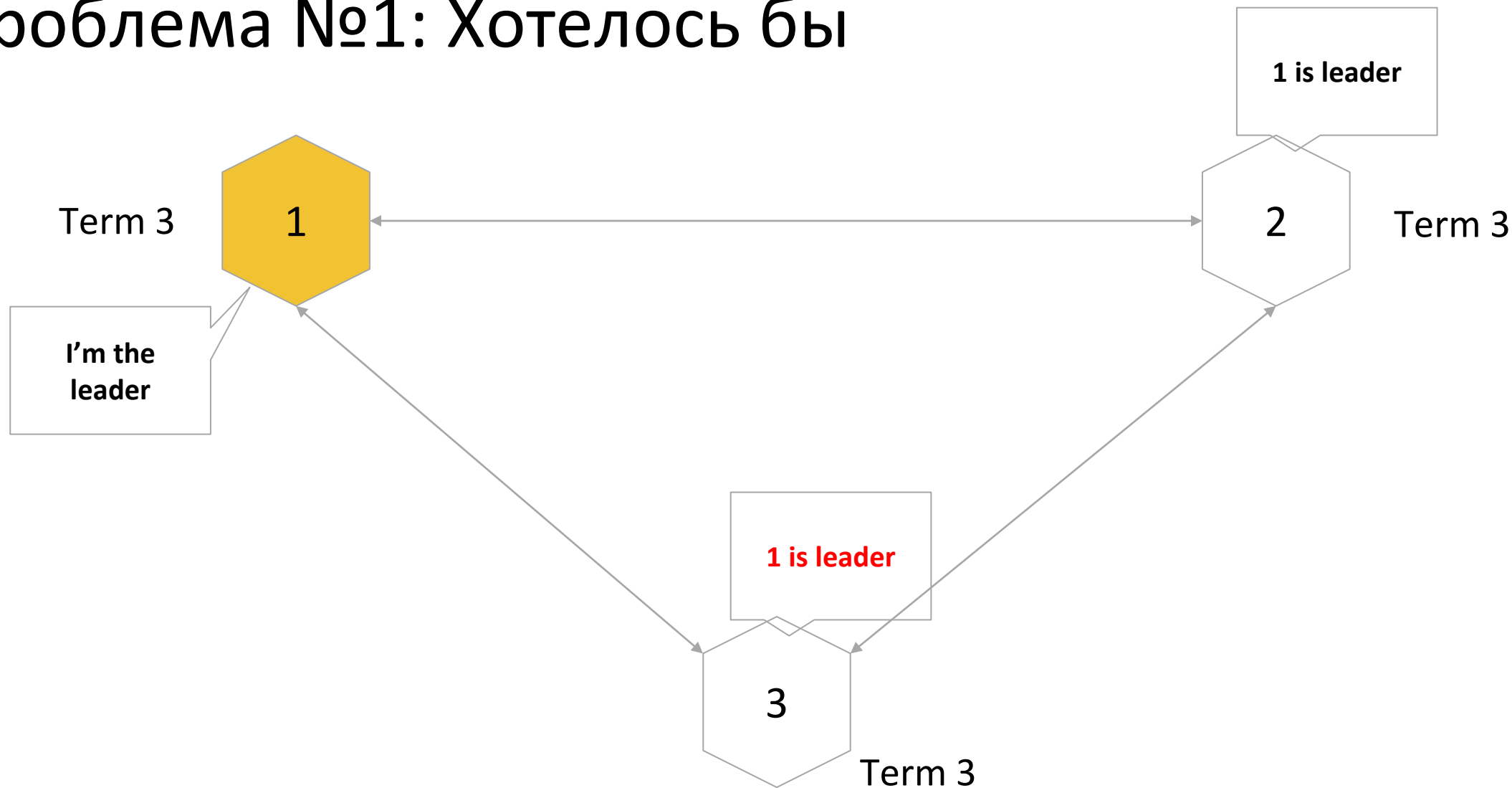
# Проблема №1: Хотелось бы



# Проблема №1: Хотелось бы

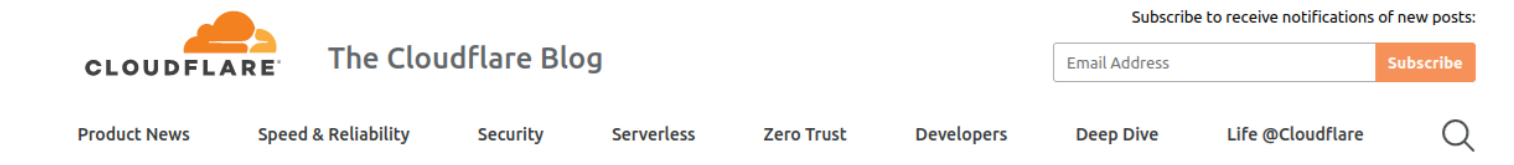


# Проблема №1: Хотелось бы



# Pre-Vote: Проблема №2

Сервер, который не видит лидера, будет постоянно его прерывать

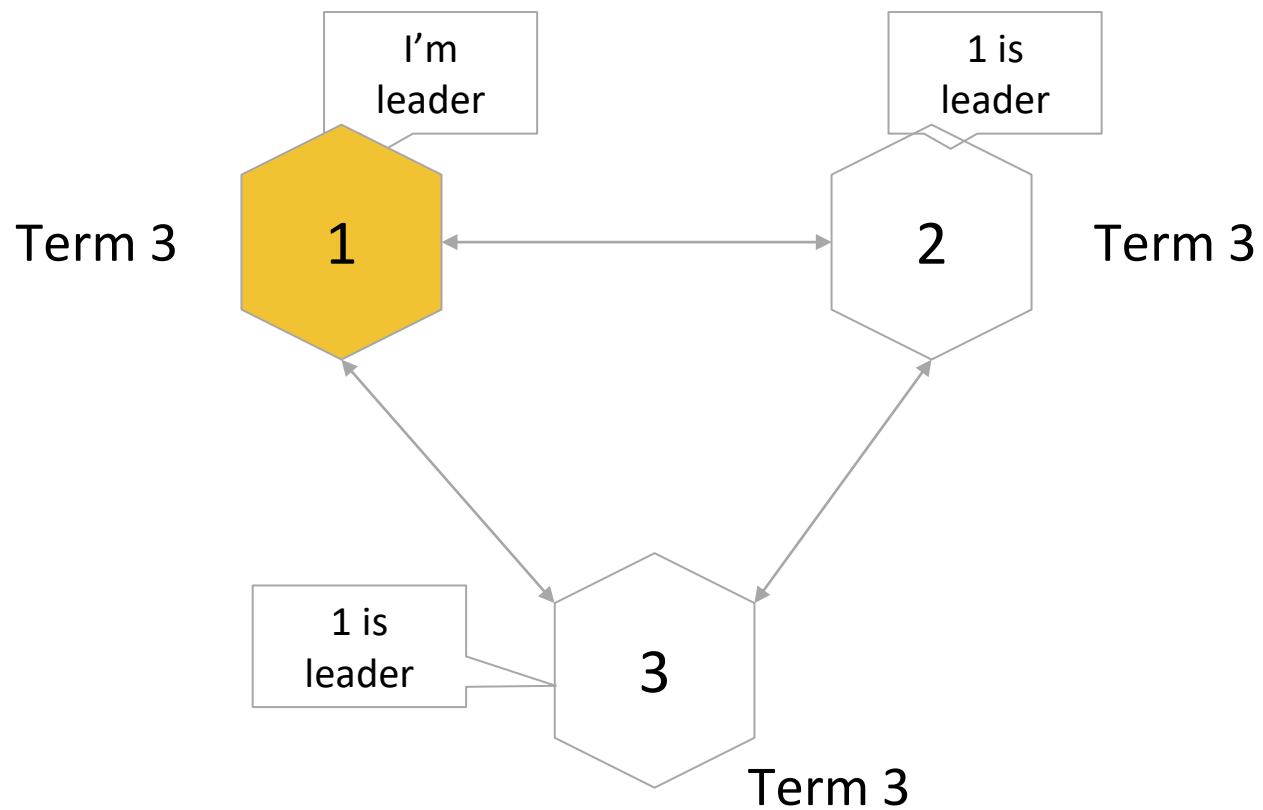


## A Byzantine failure in the real world

27.11.2020

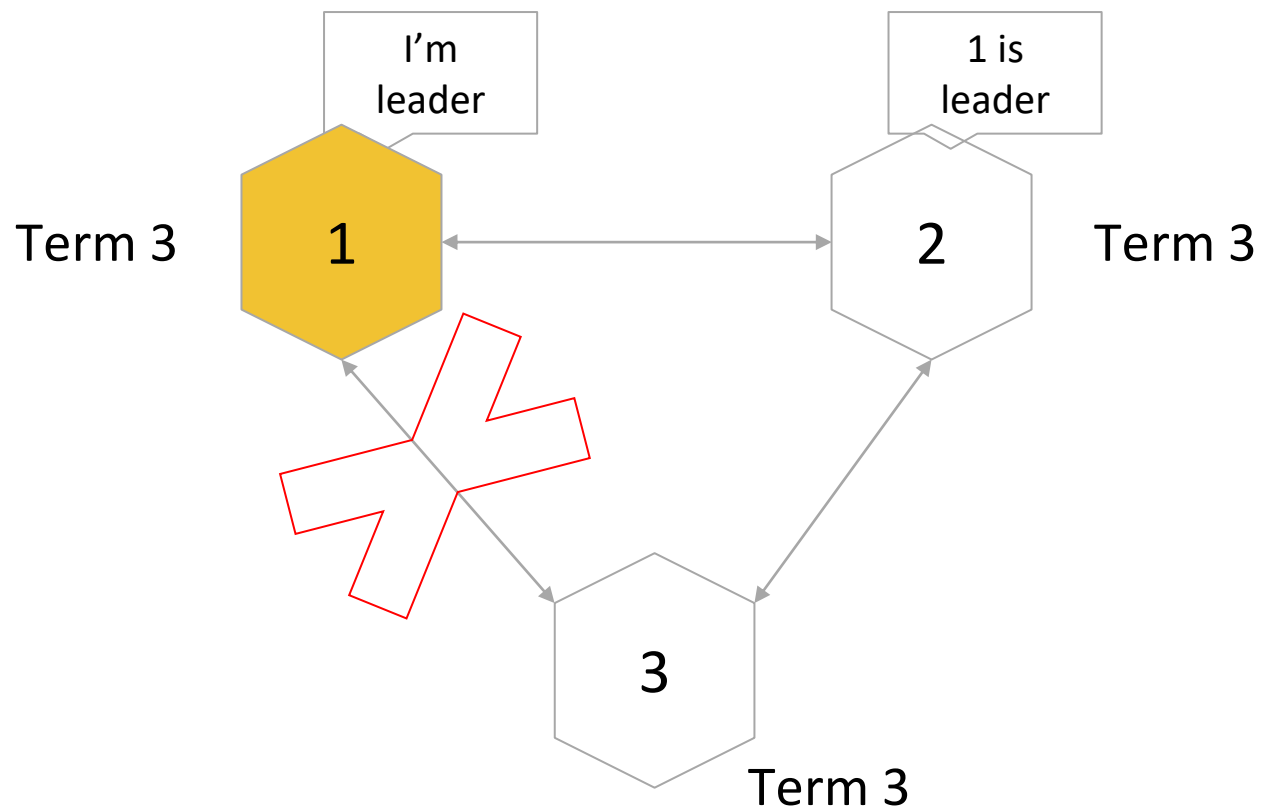
On November 2, 2020, Cloudflare had an [incident](#) that impacted the availability of the API and dashboard for six hours and 33 minutes. During this incident, the

# Проблема №2: Пример

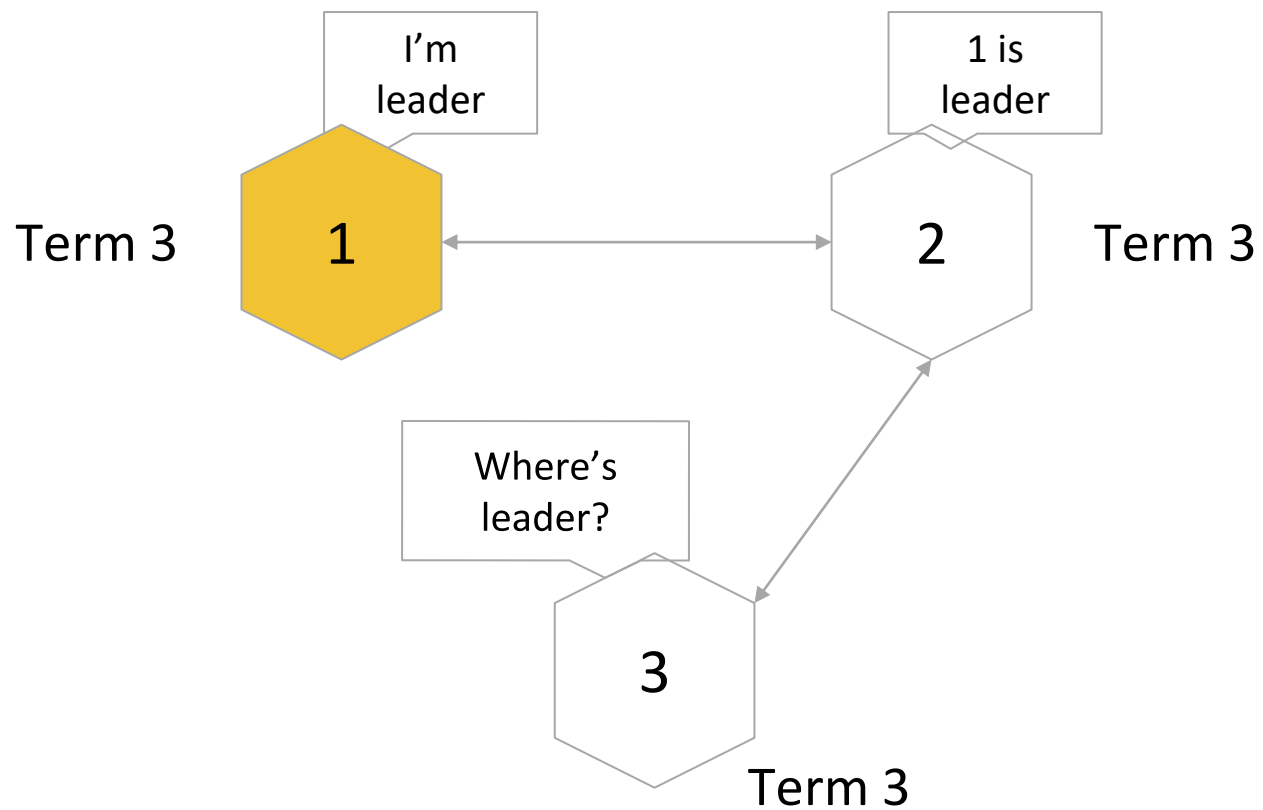




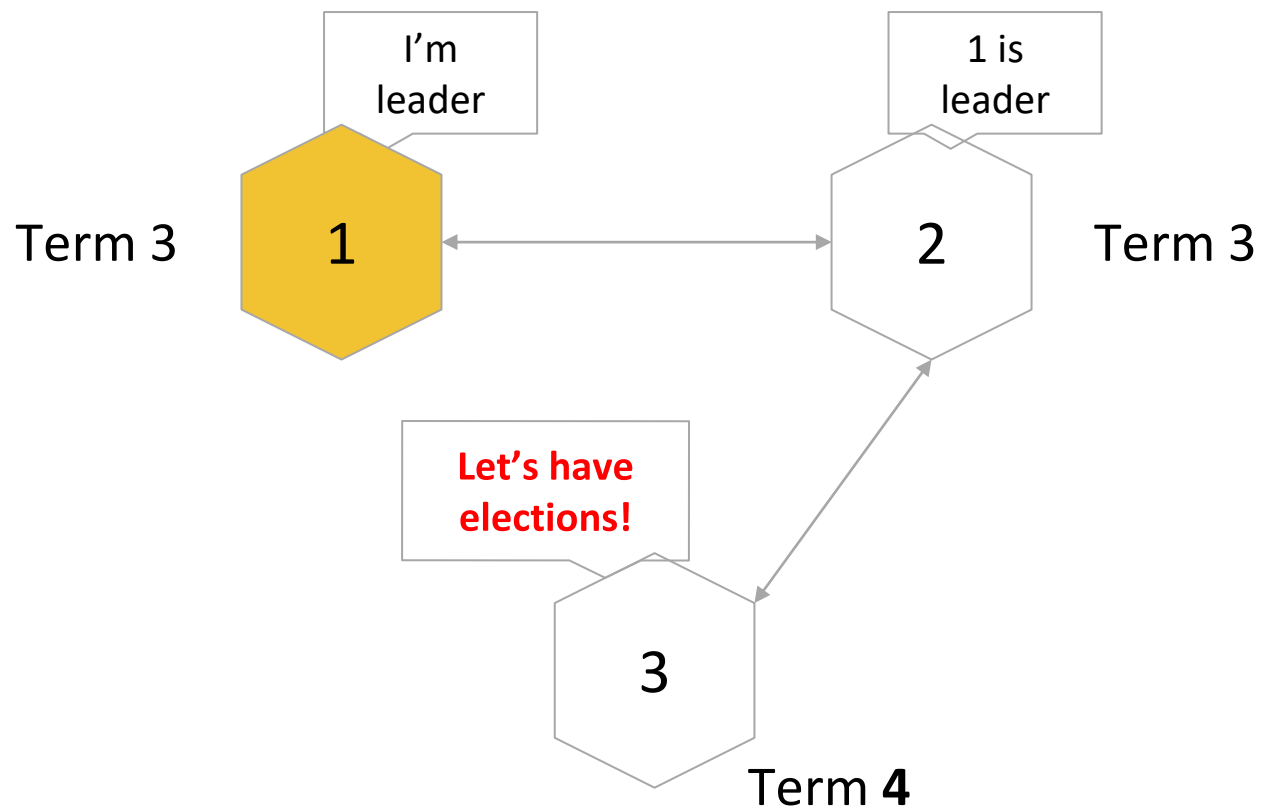
# Проблема №2: Пример



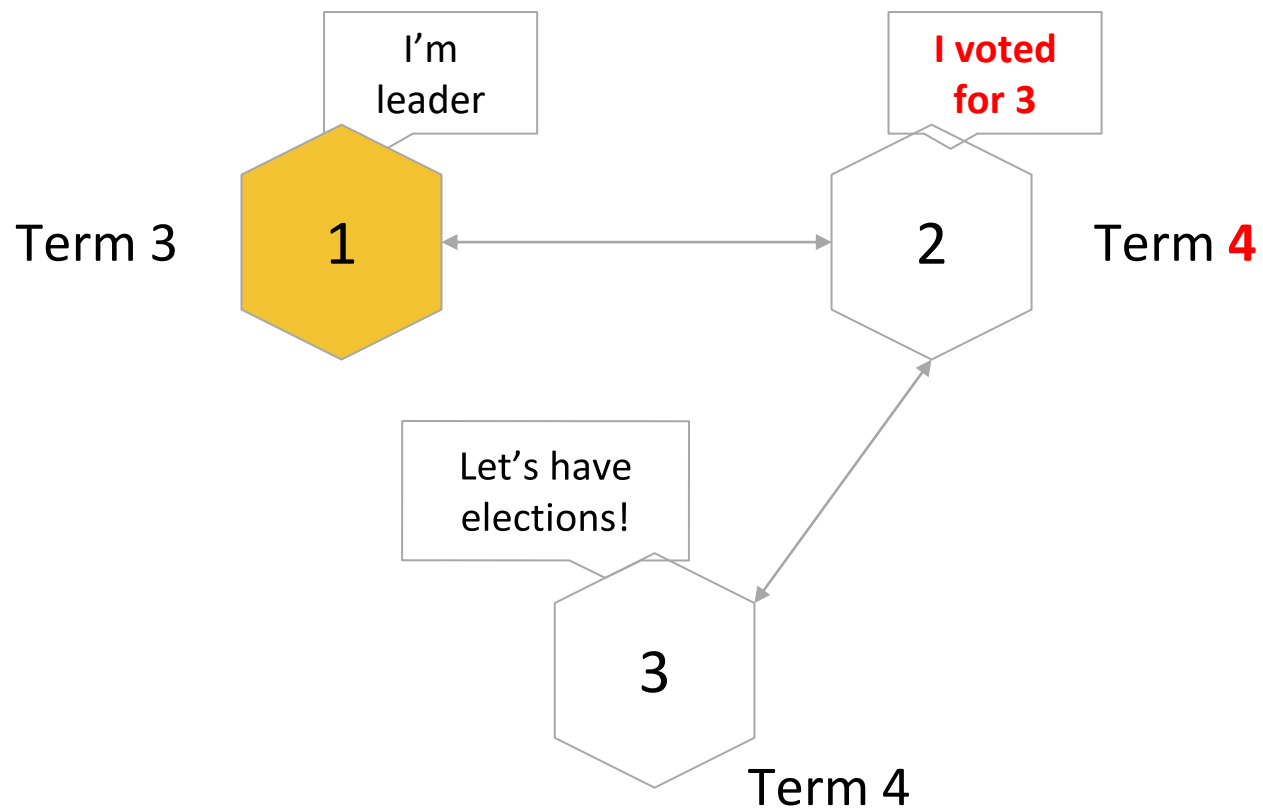
# Проблема №2: Пример



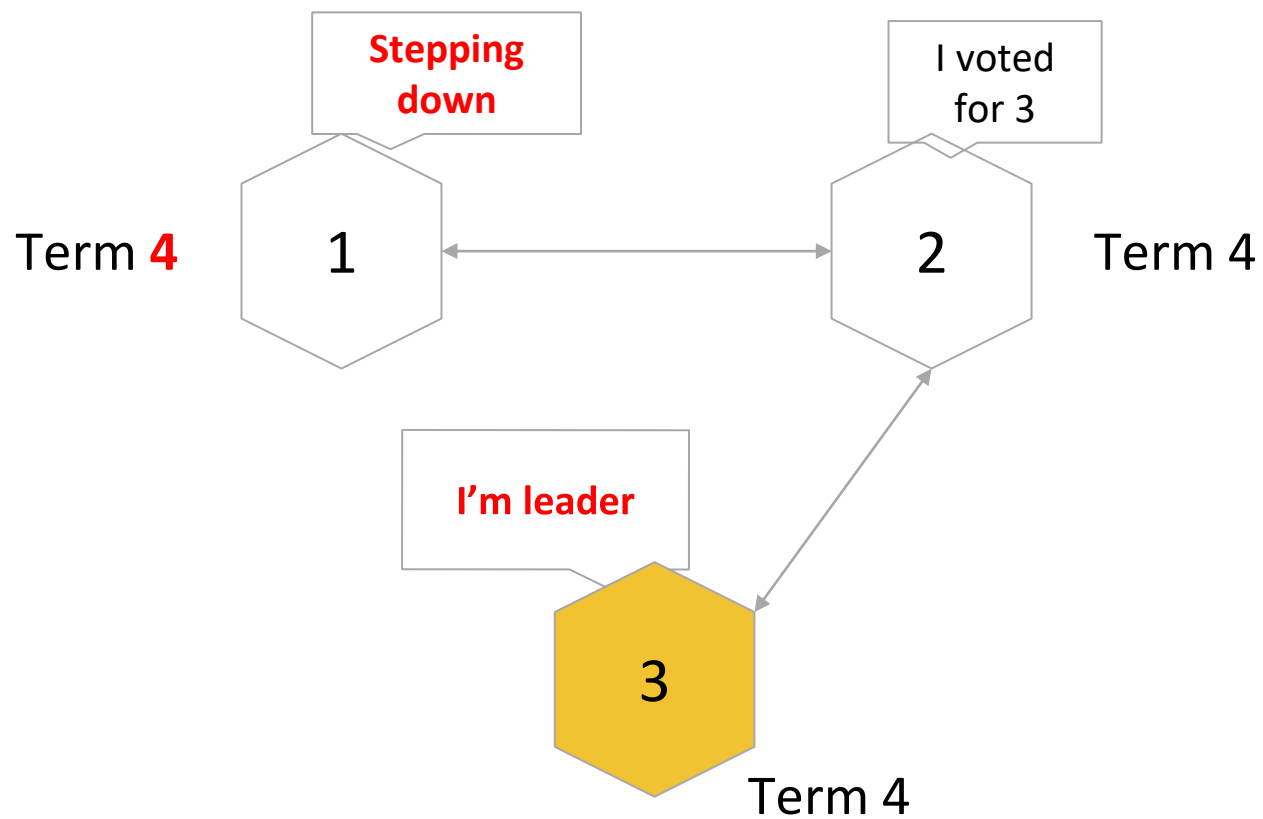
# Проблема №2: Пример



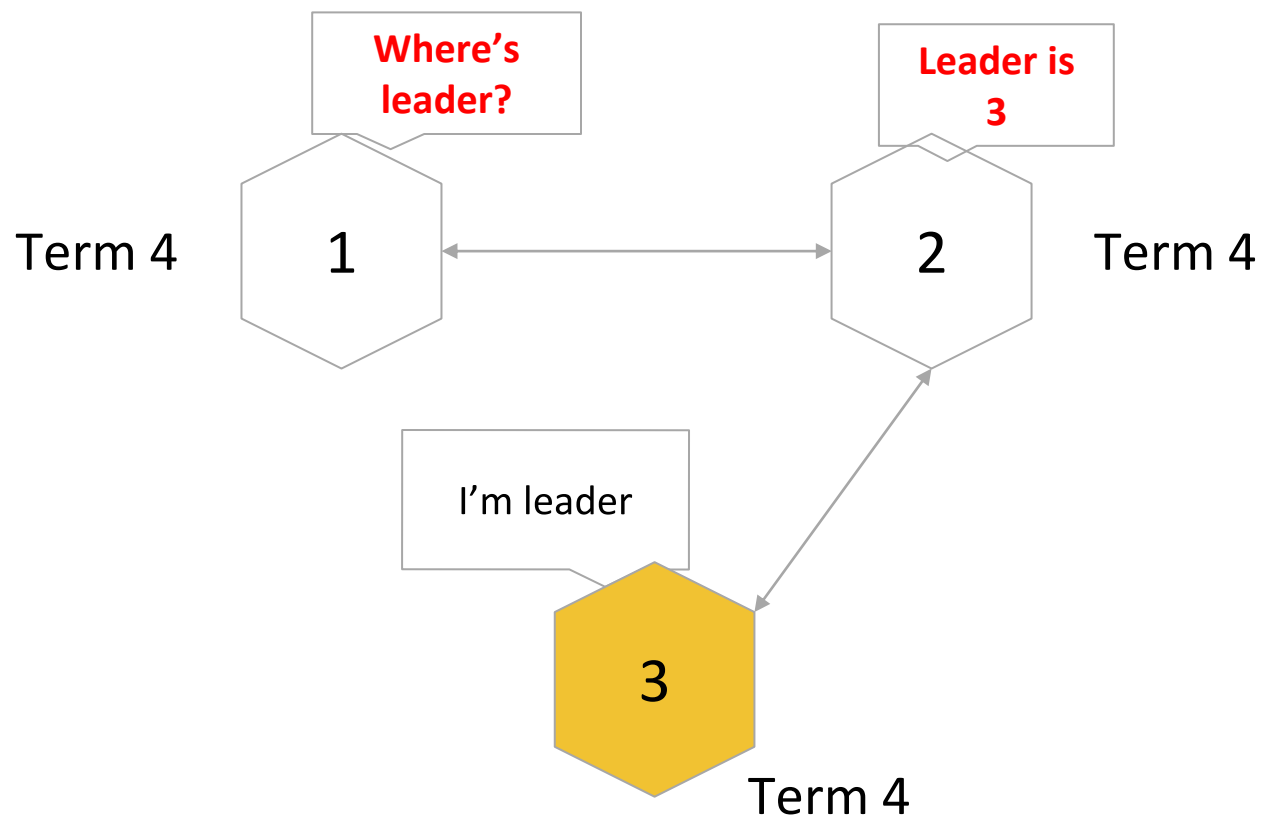
# Проблема №2: Пример



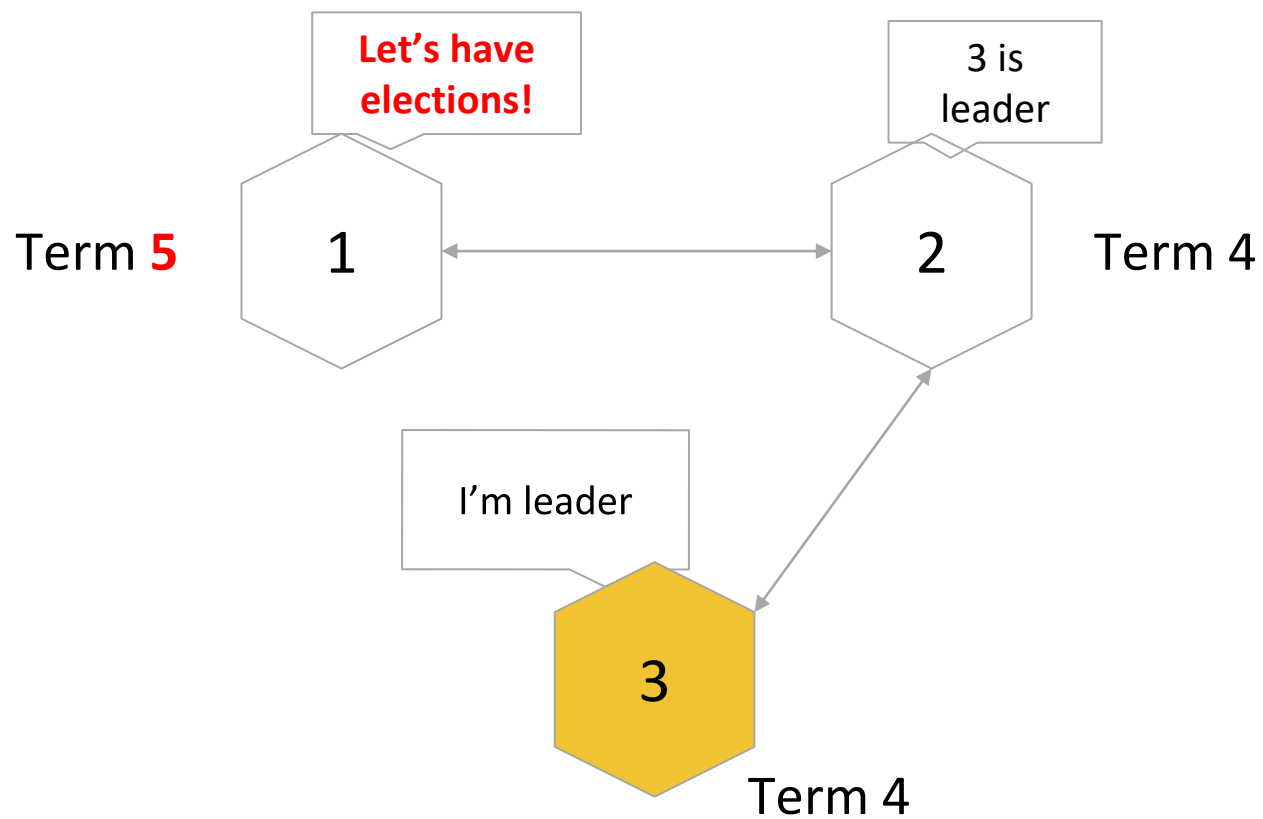
# Проблема №2: Пример



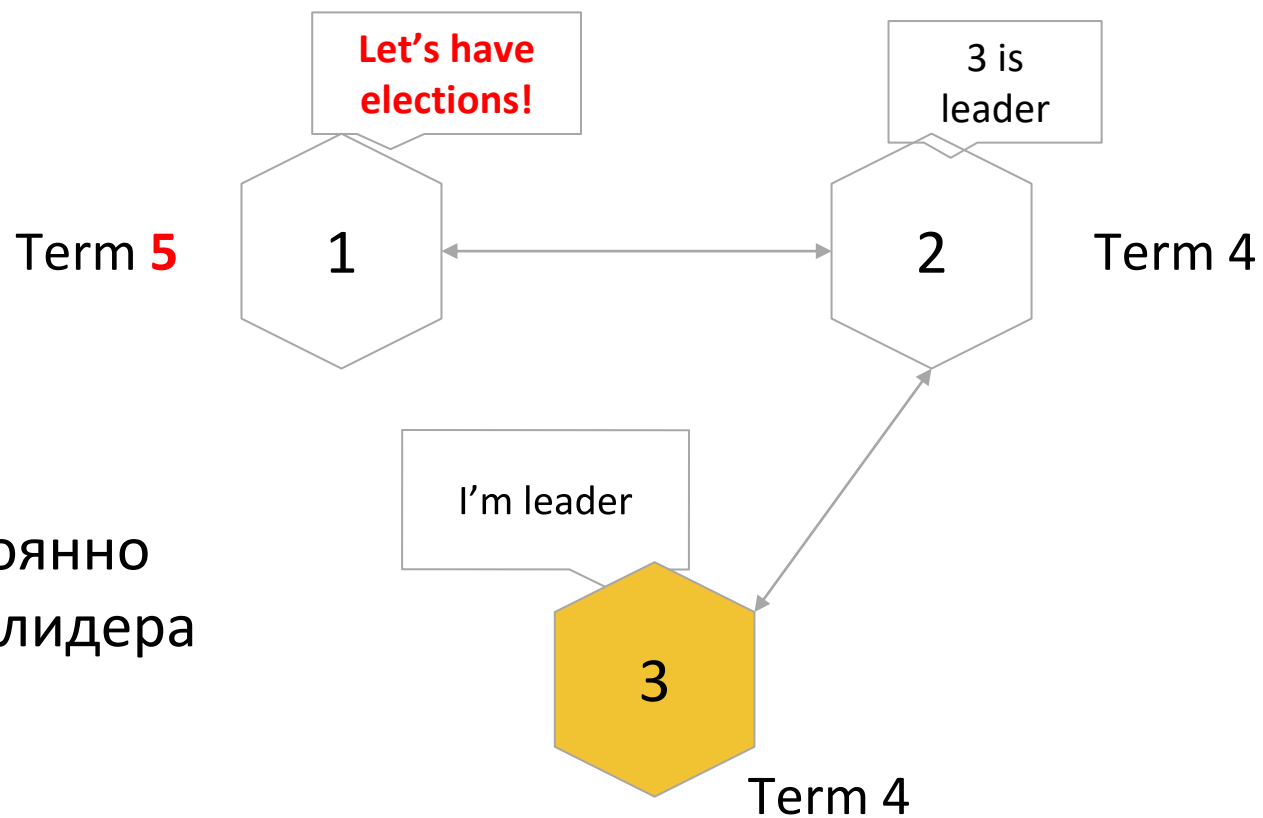
# Проблема №2: Пример



# Проблема №2: Пример



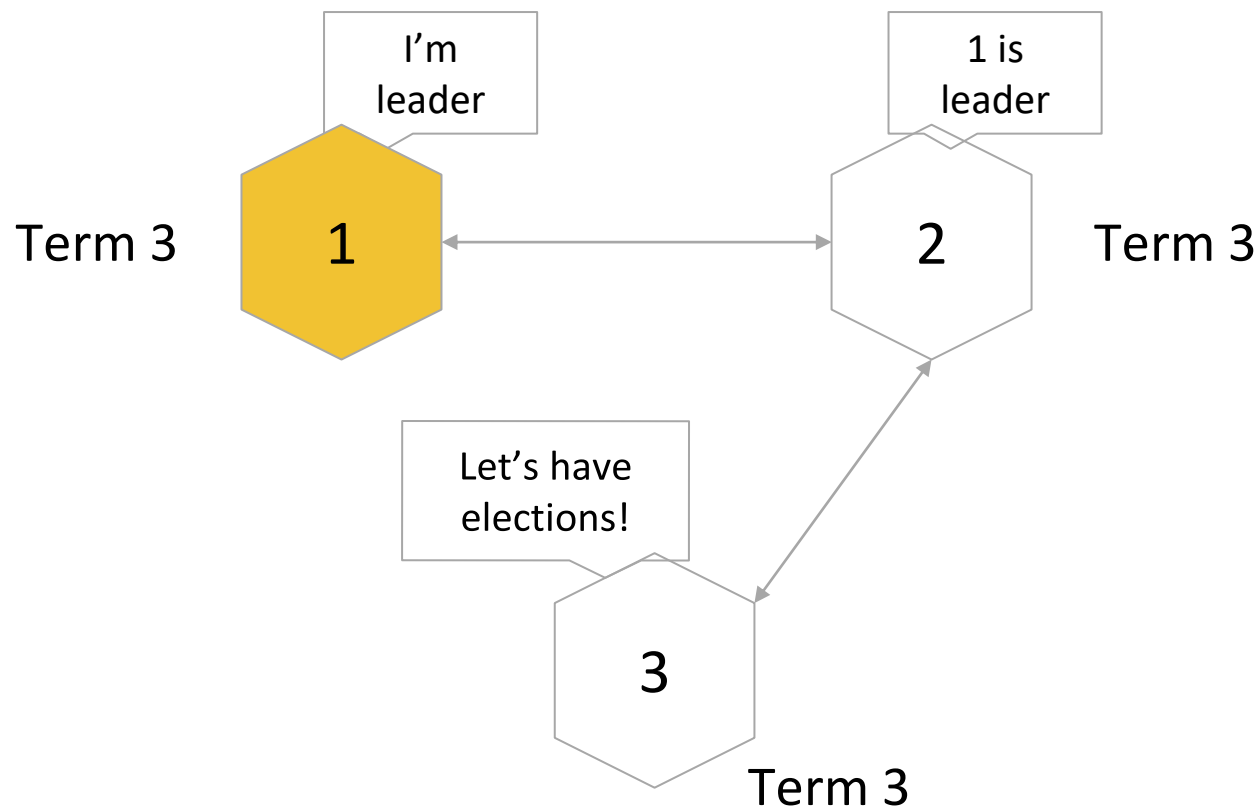
# Проблема №2: Пример



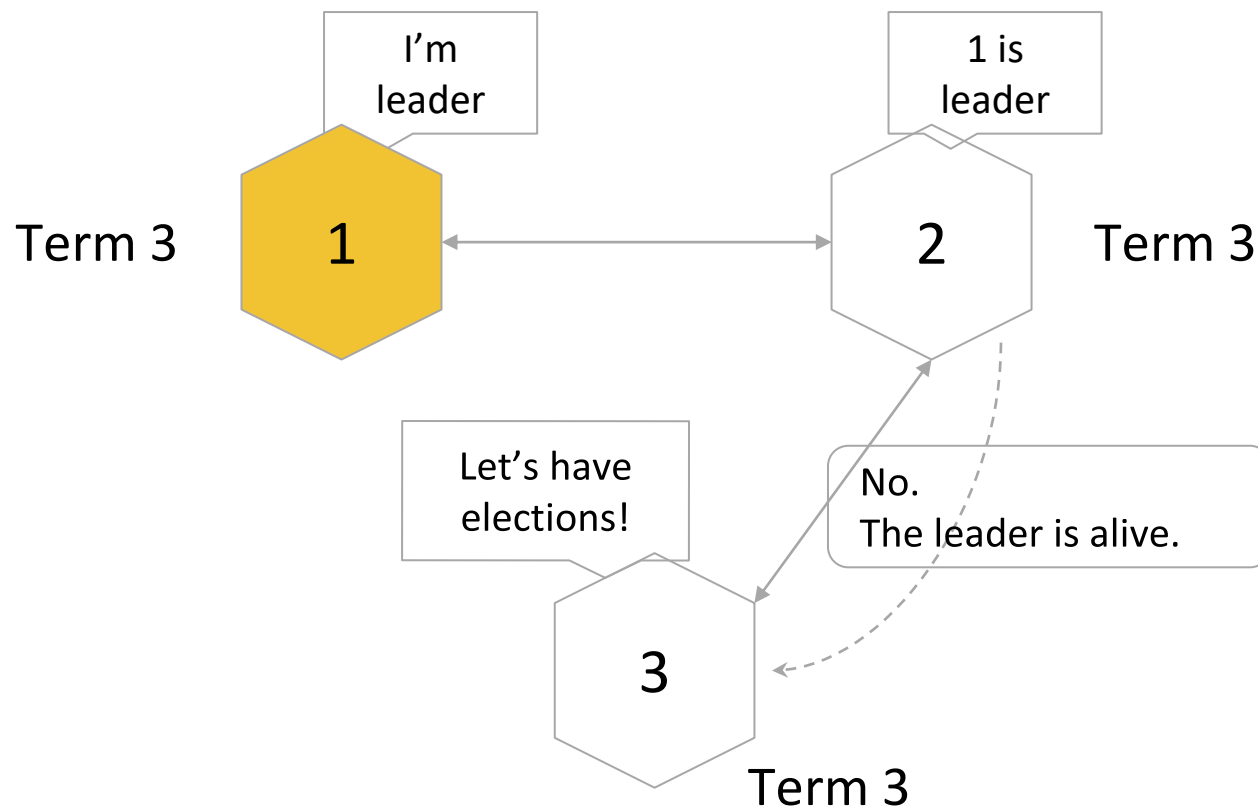
Кластер постоянно  
перевыбирает лидера



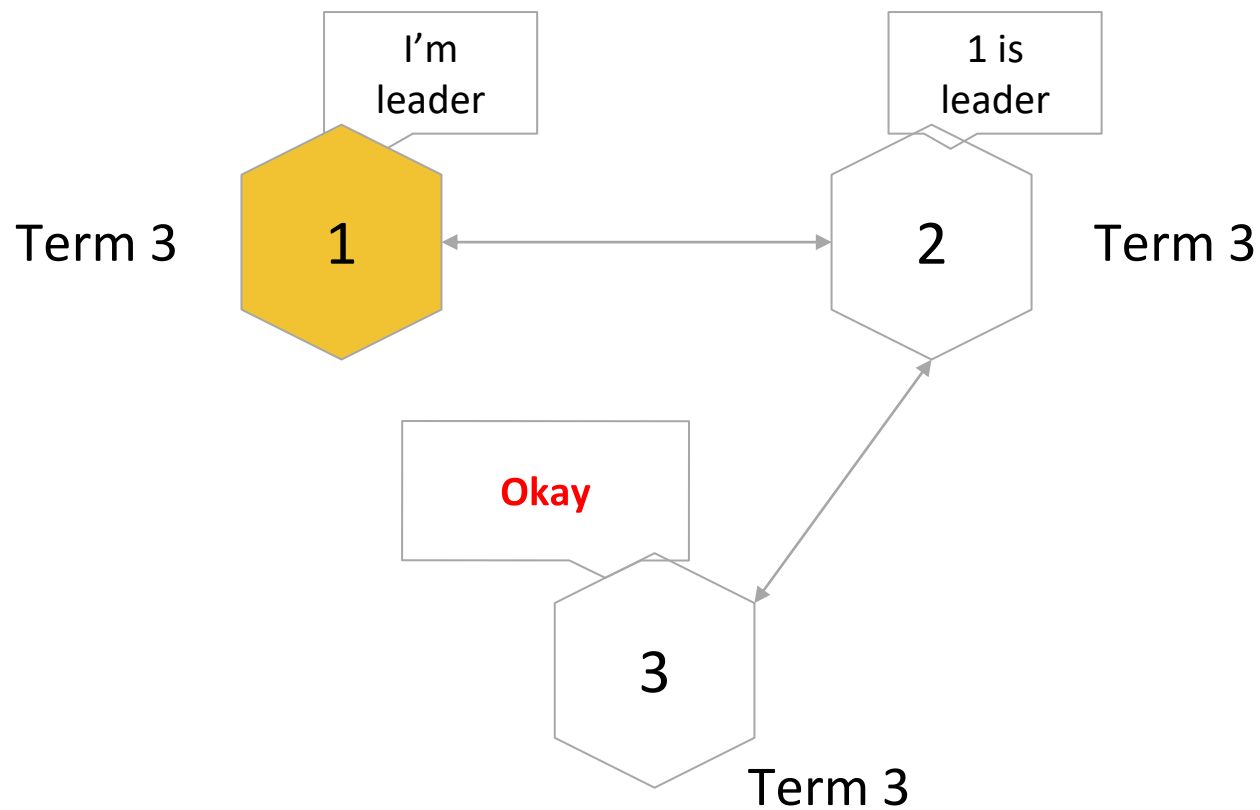
# Проблема №2: Как хочется



# Проблема №2: Как хочется



# Проблема №2: Как хочется



# Проблема №2: Как хочется

Итак, хотим, чтобы действующего лидера не мог прервать никто

# Pre-Vote: Варианты решения

## 1. Предварительные выборы

- Перед началом выборов спрашиваем, проголосуют ли за нас

2. Игнорирование выборов голосующими, которые видят лидера

3. Pre-Vote на основе метаданных

# Pre-Vote: Варианты решения

## 1. Предварительные выборы

- Перед началом выборов спрашиваем, проголосуют ли за нас
- Ответ положительный, если голосующий сам не видит лидера, и если проголосовал бы за кандидата в обычных выборах

2. Игнорирование выборов голосующими, которые видят лидера

3. Pre-Vote на основе метайнформации

# Pre-Vote: Варианты решения

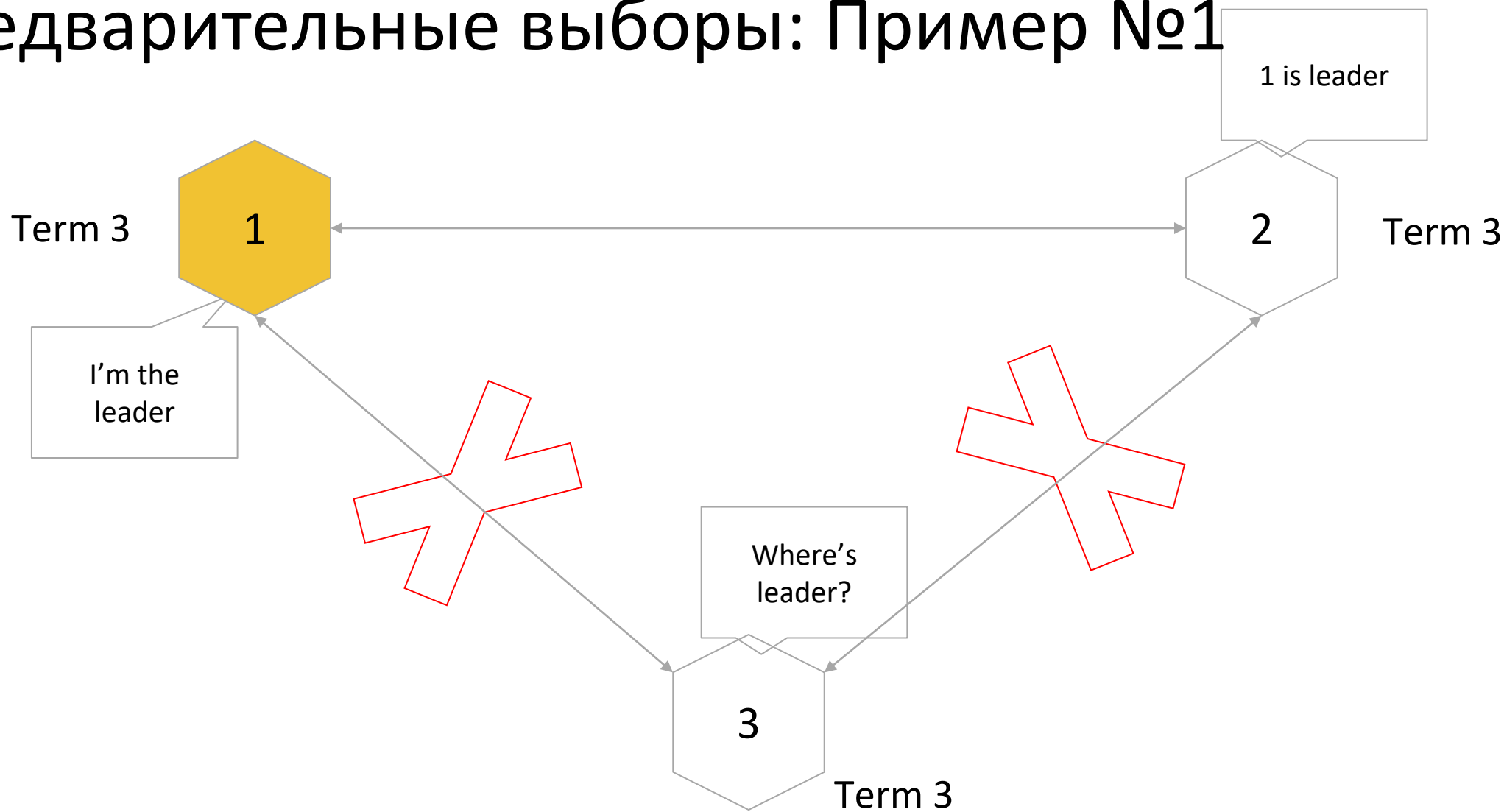
## 1. Предварительные выборы

- Перед началом выборов спрашиваем, проголосуют ли за нас
- Ответ положительный, если голосующий сам не видит лидера, и если проголосовал бы за кандидата в обычных выборах
- После получения кворума положительных ответов начинаем выборы

## 2. Игнорирование выборов голосующими, которые видят лидера

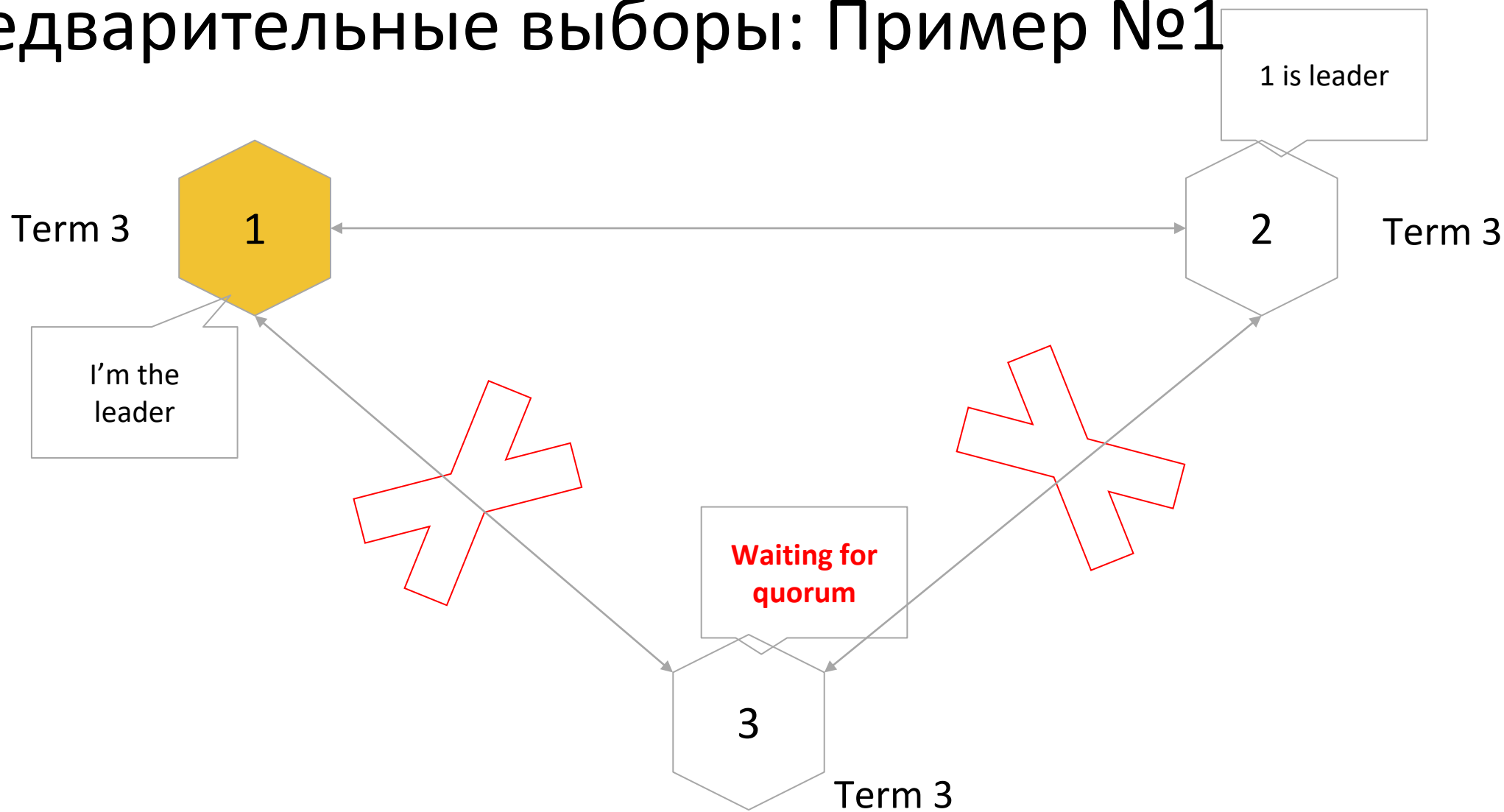
## 3. Pre-Vote на основе метайнформации

# Предварительные выборы: Пример №1

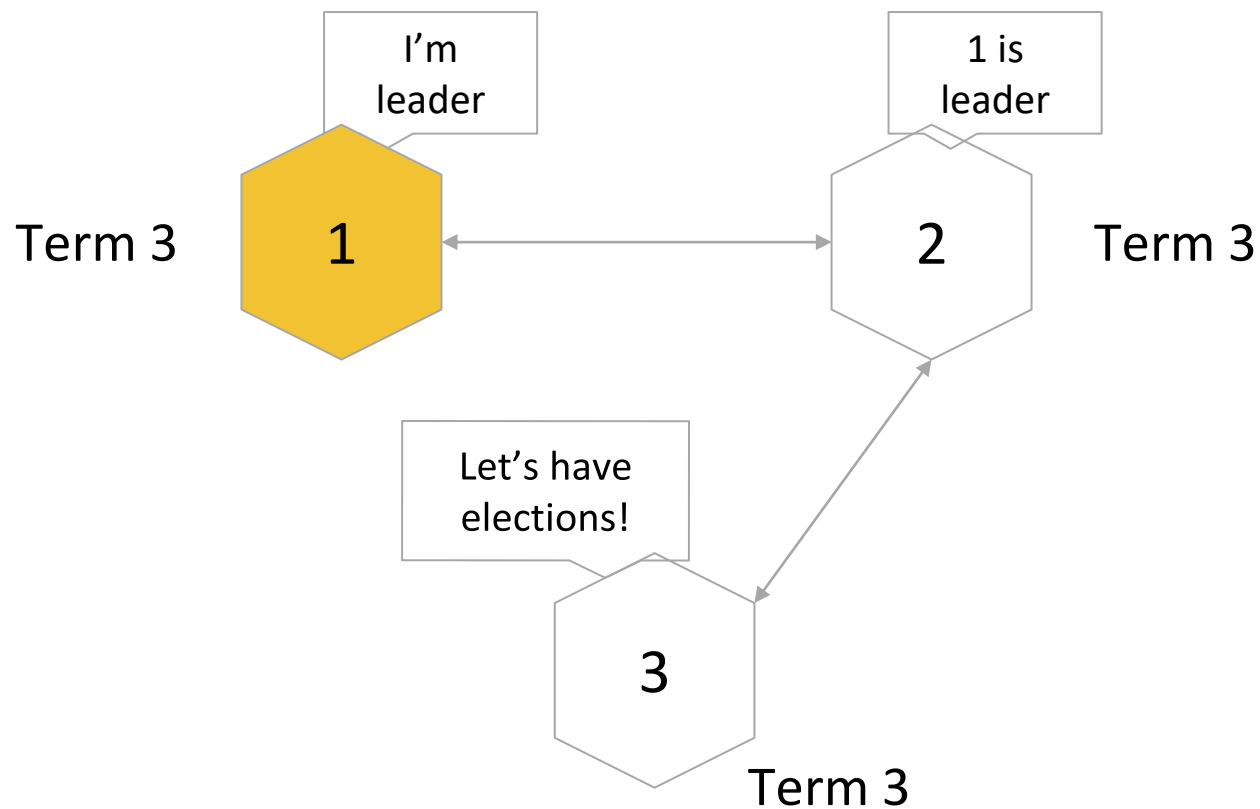




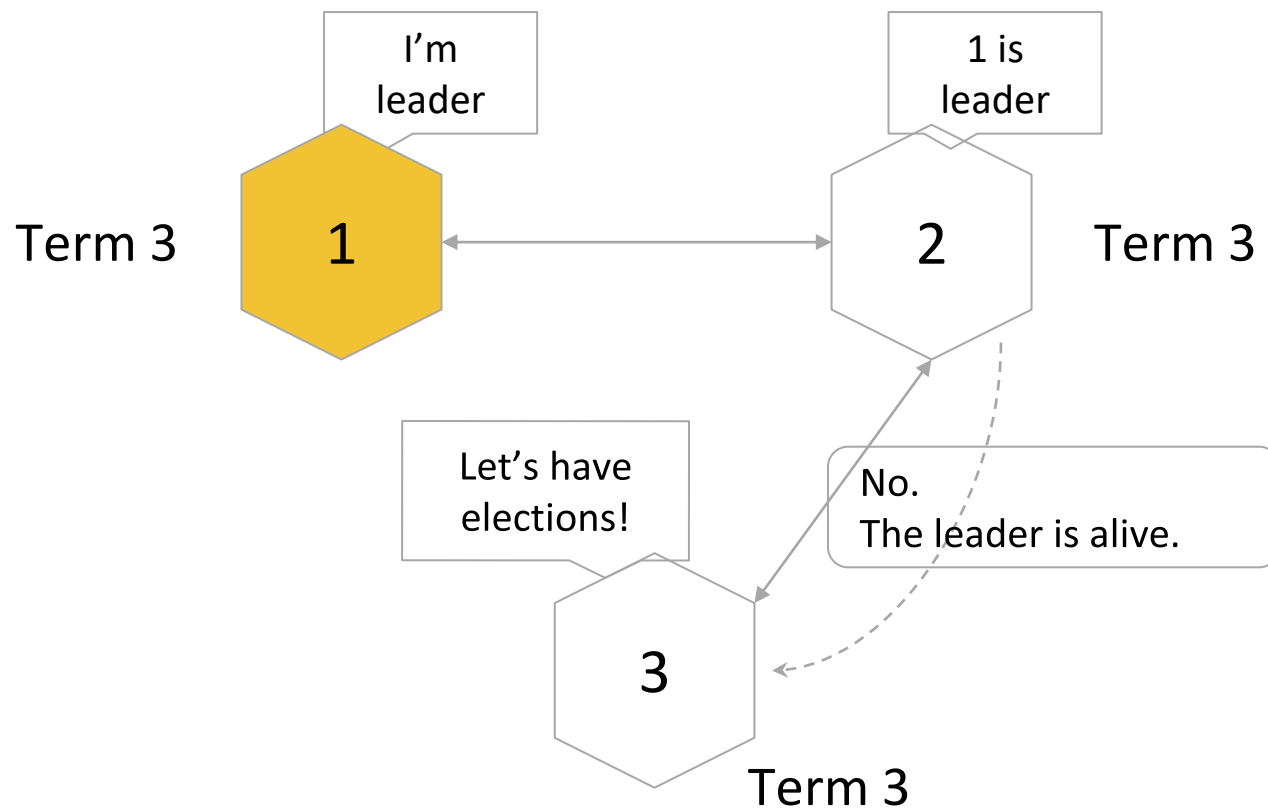
# Предварительные выборы: Пример №1



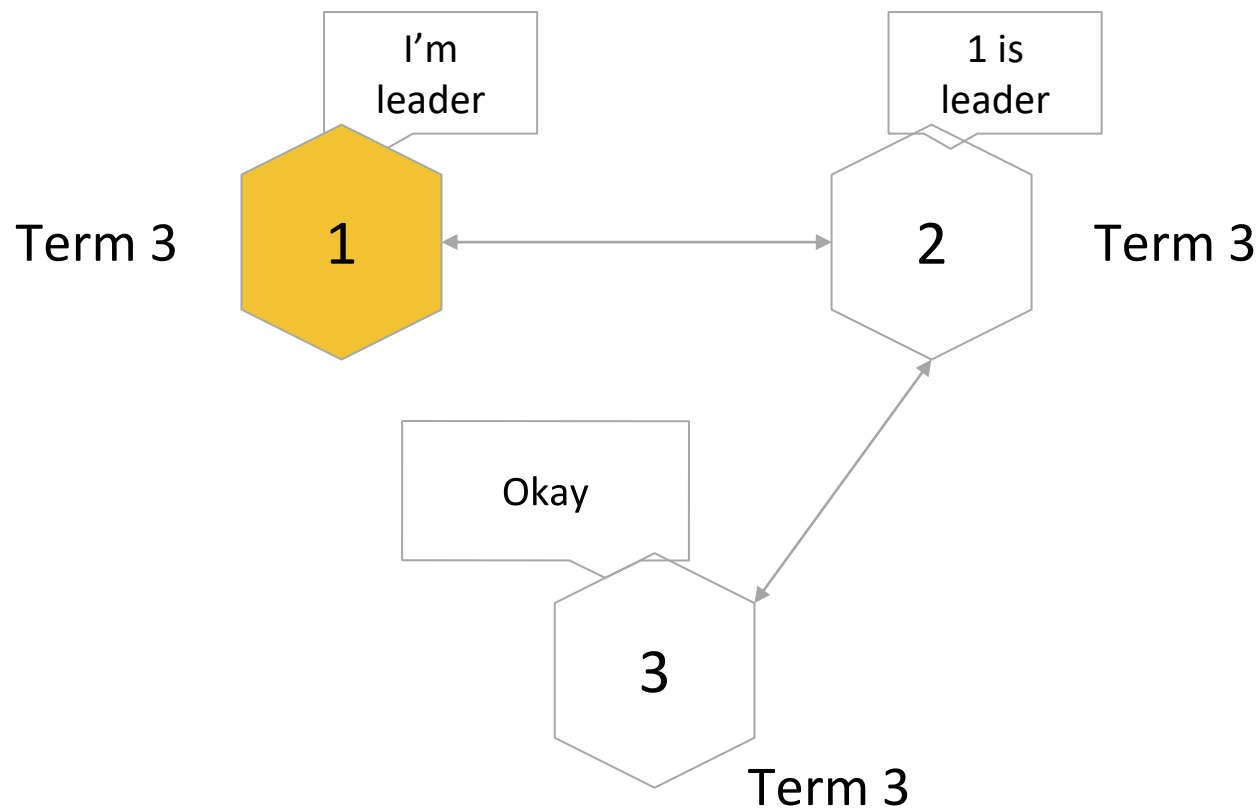
# Предварительные выборы: Пример №2



# Предварительные выборы: Пример №2



# Предварительные выборы: Пример №2



# Pre-Vote: Варианты решения

## 1. Предварительные выборы

- Перед началом выборов спрашиваем, проголосуют ли за нас
- Ответ положительный, если голосующий сам не видит лидера и если проголосовал бы за кандидата в обычных выборах
- После получения кворума положительных ответов начинаем выборы

2. Игнорирование выборов голосующими, которые видят лидера

3. Pre-Vote на основе метайнформации

Недостаток – требуется новый тип сообщений, ломающий обратную совместимость с существующими инсталляциями

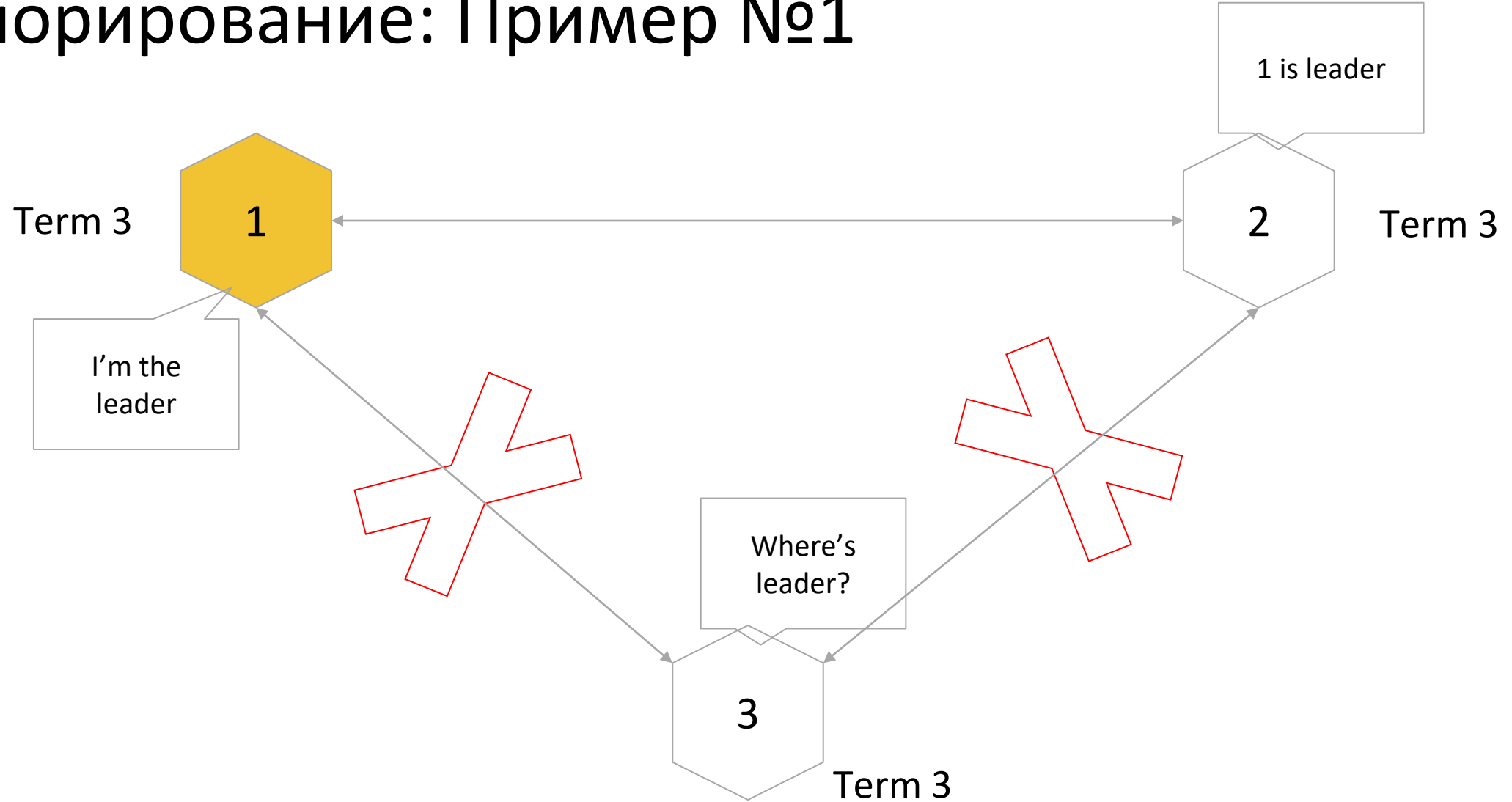
# Pre-Vote: Варианты решения

1. Предварительные выборы
2. **Игнорирование выборов голосующими, которые видят лидера**
3. Pre-Vote на основе метаданных

# Pre-Vote: Варианты решения

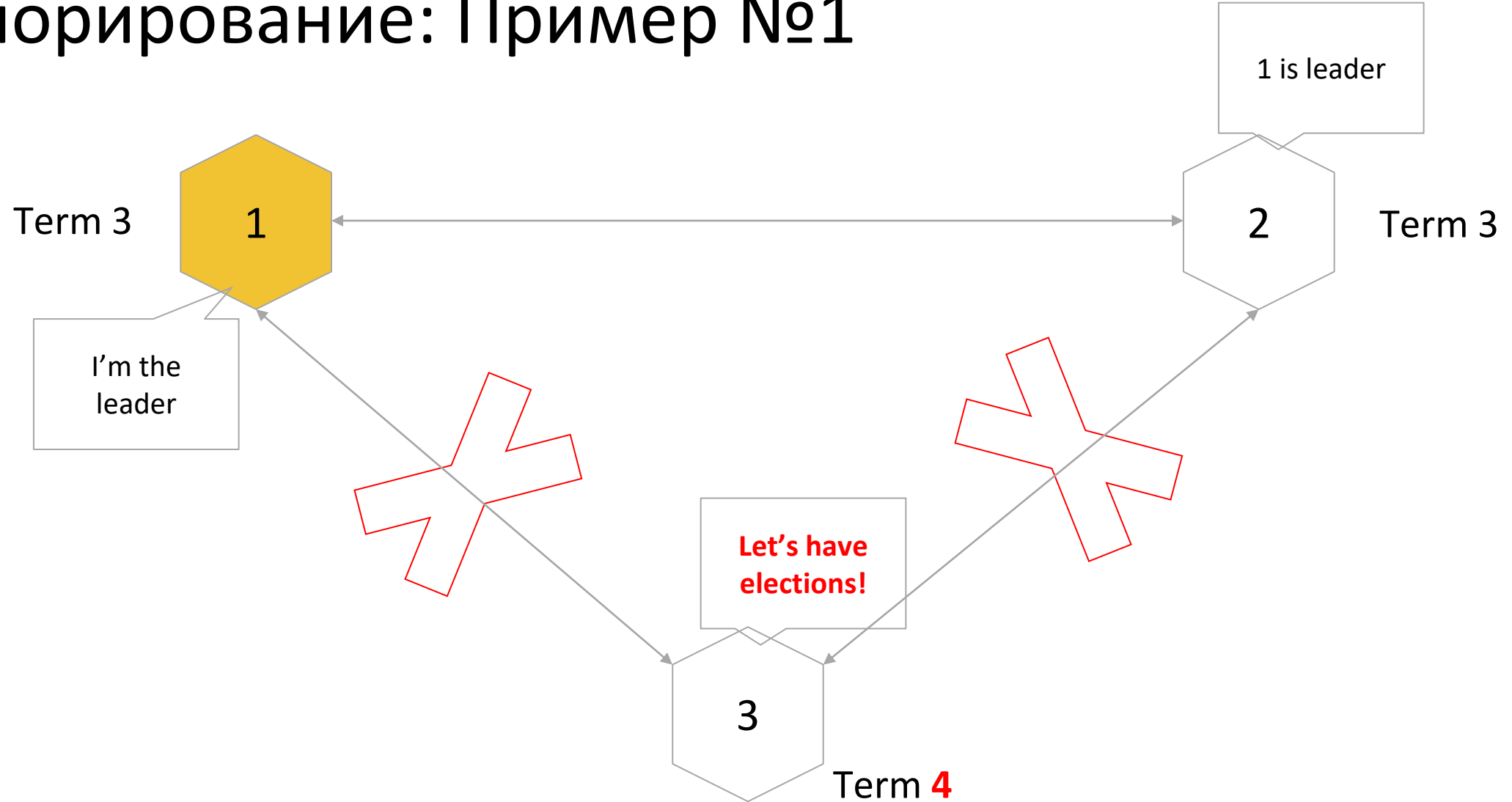
1. Предварительные выборы
2. **Игнорирование выборов голосующими, которые видят лидера**
  - Сервер, видящий лидера напрямую, игнорирует любые запросы на голос от соседей
3. Pre-Vote на основе метаданных

# Игнорирование: Пример №1

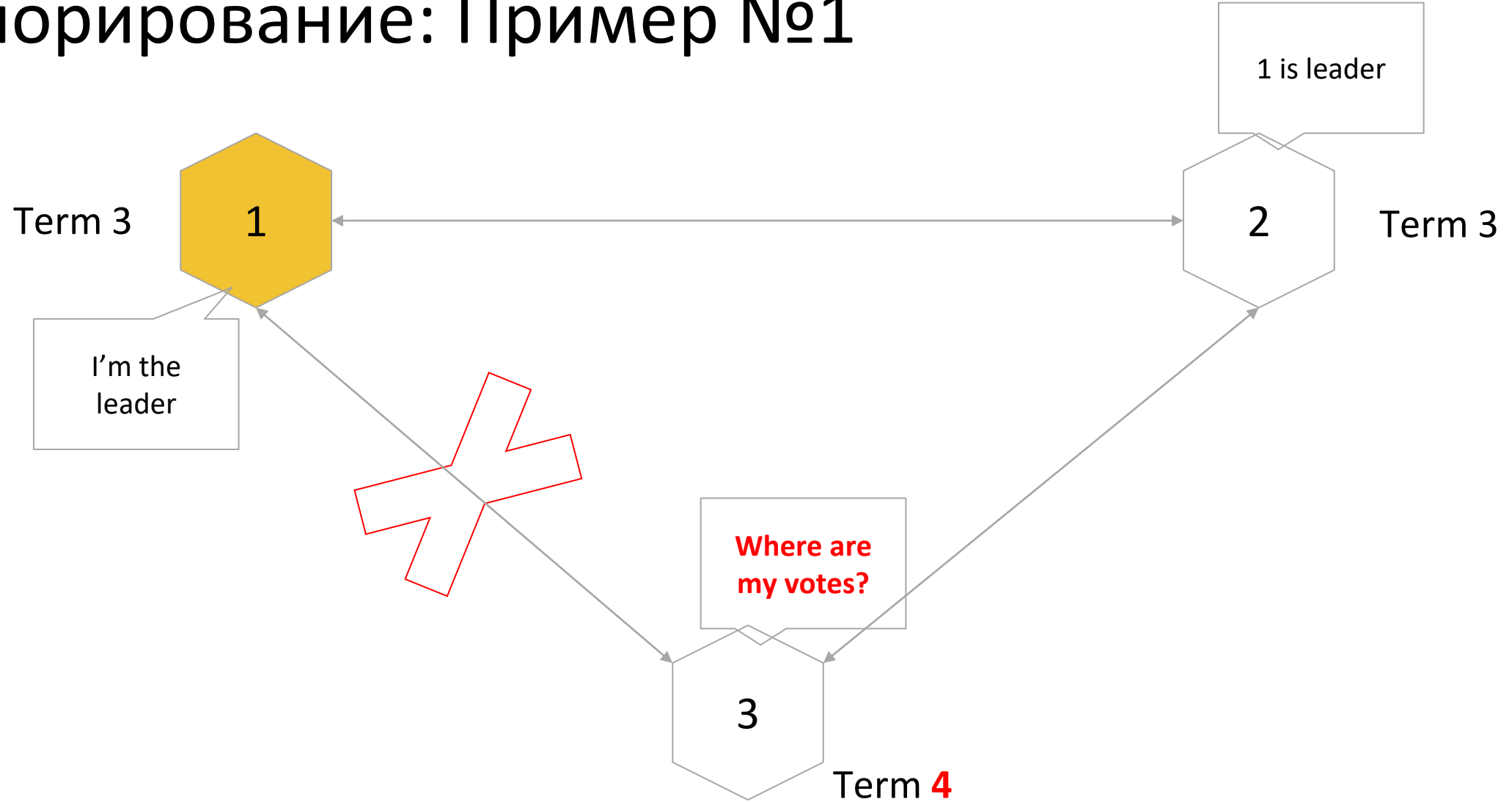




# Игнорирование: Пример №1



# Игнорирование: Пример №1



# Игнорирование: Пример №1



# Pre-Vote: Варианты решения

1. Предварительные выборы
2. **Игнорирование выборов голосующими, которые видят лидера**
  - Сервер, видящий лидера напрямую, игнорирует любые запросы на голос от соседей
3. Pre-Vote на основе метаданных

Недостаток – решение неполное!

# Pre-Vote: Варианты решения

1. Предварительные выборы
2. Игнорирование выборов голосующими, которые видят лидера
3. **Pre-Vote на основе метаданных**
  - Каждый сервер знает состояние Raft соседей за счет широковещательной рассылки

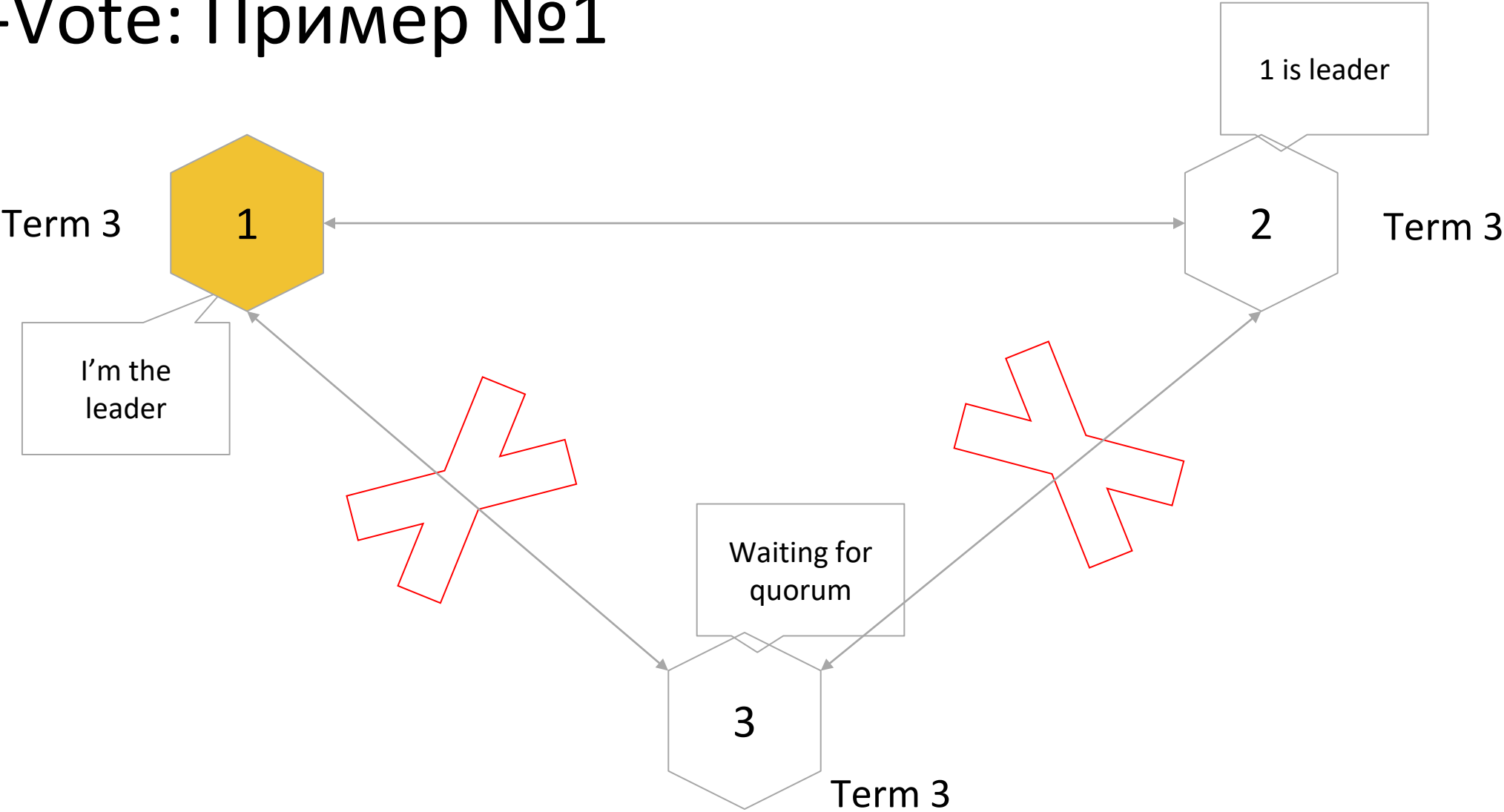
# Pre-Vote: Варианты решения

1. Предварительные выборы
2. Игнорирование выборов голосующими, которые видят лидера
3. **Pre-Vote на основе метаданных**
  - Каждый сервер знает состояние Raft соседей за счет широковещательной рассылки
  - Добавим в сообщение о состоянии информацию о том, виден ли лидер

# Pre-Vote: Варианты решения

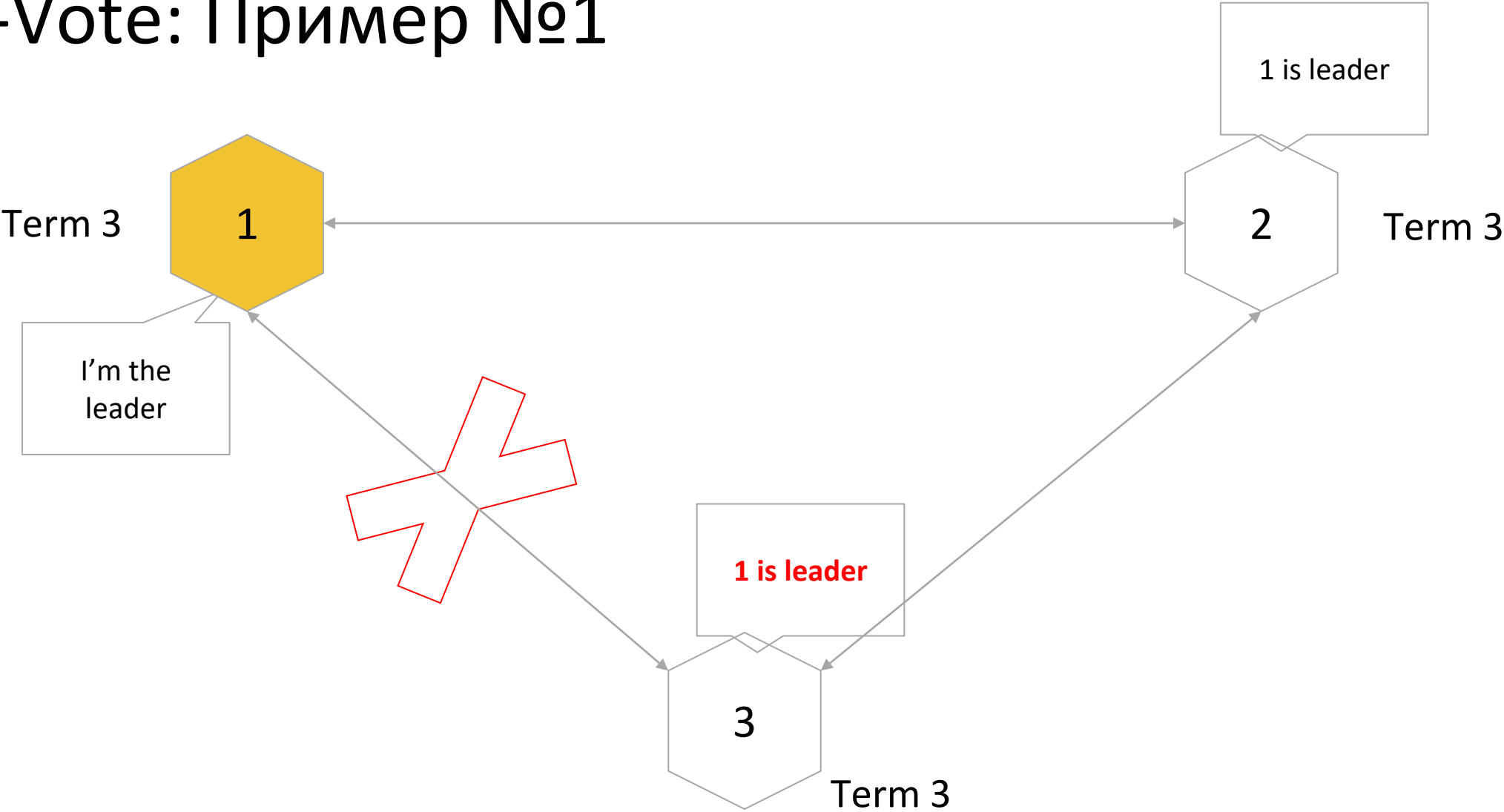
1. Предварительные выборы
2. Игнорирование выборов голосующими, которые видят лидера
3. **Pre-Vote на основе метаинформации**
  - Каждый сервер знает состояние Raft соседей за счет широковещательной рассылки
  - Добавим в сообщение о состоянии информацию о том, виден ли лидер
  - Если кто-то видит лидера в текущем терме – не начинаем выборы

# Pre-Vote: Пример №1

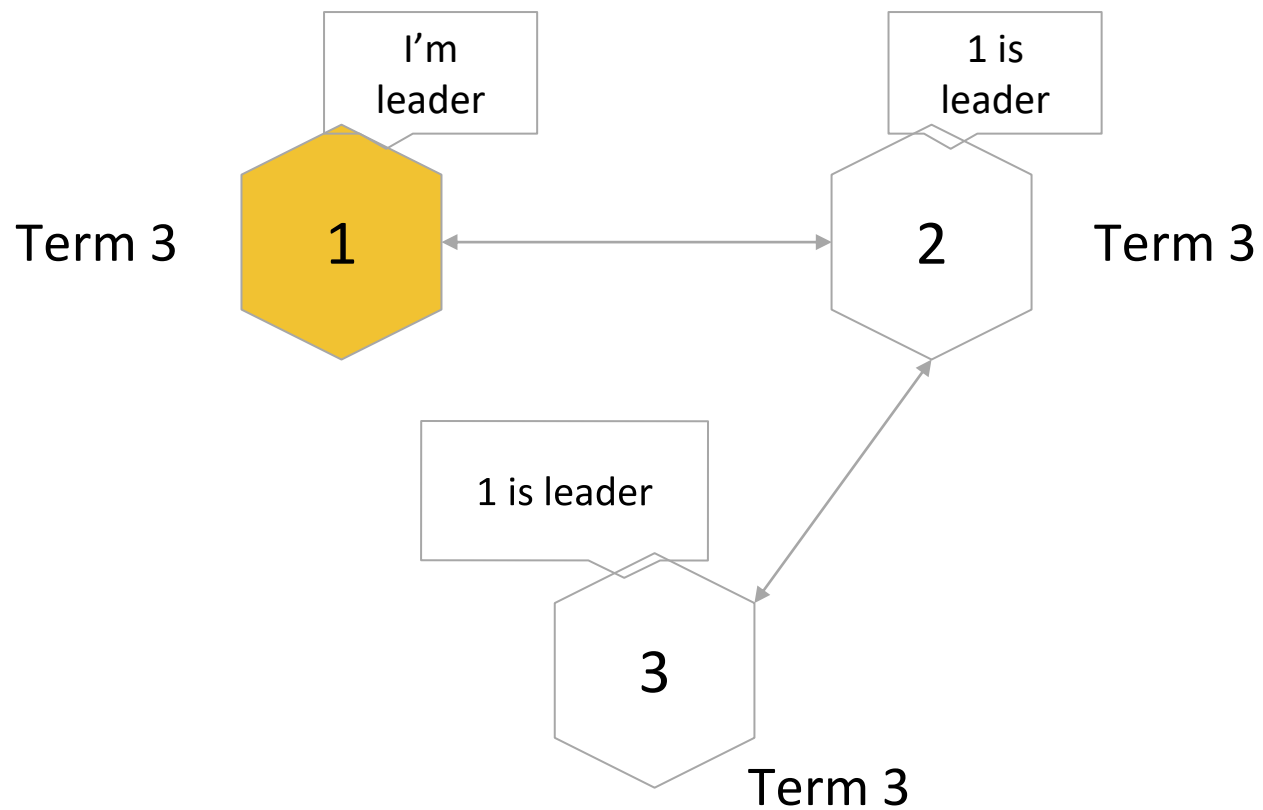




# Pre-Vote: Пример №1



# Pre-Vote: Пример №2



# Pre-Vote: Варианты решения

1. Предварительные выборы
2. Игнорирование выборов голосующими, которые видят лидера
3. **Pre-Vote на основе метаинформации**
  - Каждый сервер знает состояние Raft соседей за счет широковещательной рассылки
  - Добавим в сообщение о состоянии информацию о том, виден ли лидер
  - Если кто-то видит лидера в текущем терме - не начинаем выборы

Решение как с предварительными выборами, но обратно совместимое

# Меню

- ✓ Raft overview
  - Термины: Journal, Term, LSN
  - Выборы
  - Гарантии
  - Ожидания != Реальность
- ✓ Raft / Tarantool: особенности
  - Настройки Raft
    - ✓ Pre-Vote
      - **Split-Vote detection**
      - Fencing

# Split-Vote: Проблема

Долгая недоступность кластера на запись после потери лидера

# Split-Vote: Проблема

Долгая недоступность кластера на запись после потери лидера

Каждый раунд выборов, в результате которого не выбран лидер –  $\geq 5$  секунд простоя

# Split-Vote: Проблема

Долгая недоступность кластера на запись после потери лидера

Каждый раунд выборов, в результате которого не выбран лидер –  $\geq 5$  секунд простоя

В случае успешных выборов лидер появляется в самом начале раунда

# Split-Vote: Проблема

Долгая недоступность кластера на запись после потери лидера

Каждый раунд выборов, в результате которого не выбран лидер –  $\geq 5$  секунд простоя

В случае успешных выборов лидер появляется в самом начале раунда

Если голоса разделились, оставшееся время раунда – простой



# Split-Vote: Решение

У нас есть широковещательная рассылка отданных голосов

# Split-Vote: Решение

У нас есть широковещательная рассылка отданных голосов

Каждый из узлов может вести подсчёт голосов за всех кандидатов

# Split-Vote: Решение

У нас есть широковещательная рассылка отданных голосов

Каждый из узлов может вести подсчёт голосов за всех кандидатов

Стало понятно, что выборы никто не выиграет? Новый раунд

# Split-Vote: Решение

У нас есть широковещательная рассылка отданных голосов

Каждый из узлов может вести подсчёт голосов за всех кандидатов

Стало понятно, что выборы никто не выиграет? Новый раунд

Экономия до 90% времени на каждый раунд с разделившимися голосами

# Меню

- ✓ Raft overview
  - Термины: Journal, Term, LSN
  - Выборы
  - Гарантии
  - Ожидания != Реальность
- ✓ Raft / Tarantool: особенности
  - Настройки Raft
    - ✓ Pre-Vote
    - ✓ Split-Vote detection
      - **Fencing**

# Fencing: Проблема

- Для записи в кластер необходимо обратиться к текущему лидеру

# Fencing: Проблема

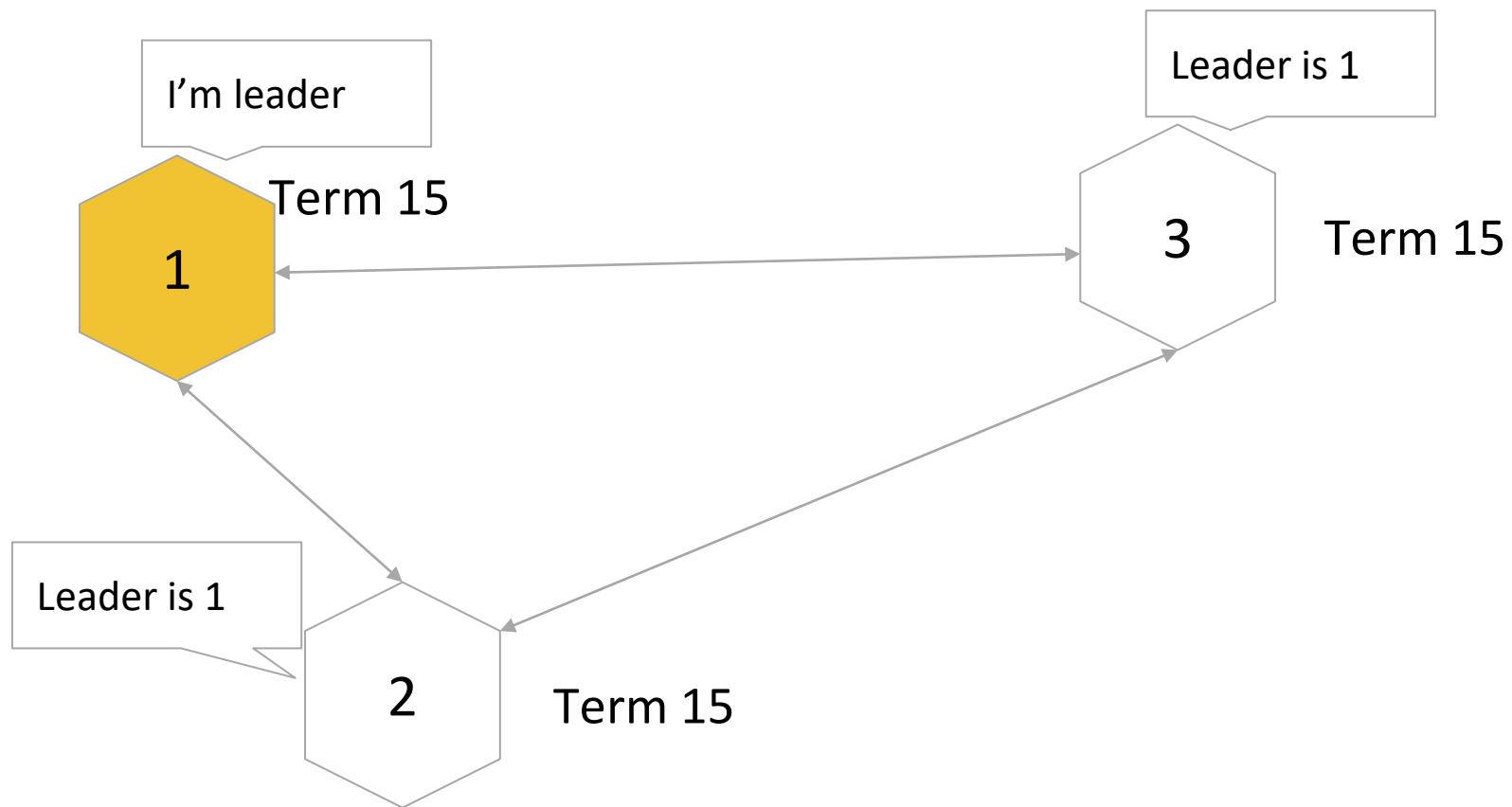
- Для записи в кластер необходимо обратиться к текущему лидеру
- Возможность существования нескольких нод, считающих себя лидерами одновременно (в разных term)

# Fencing: Проблема №1

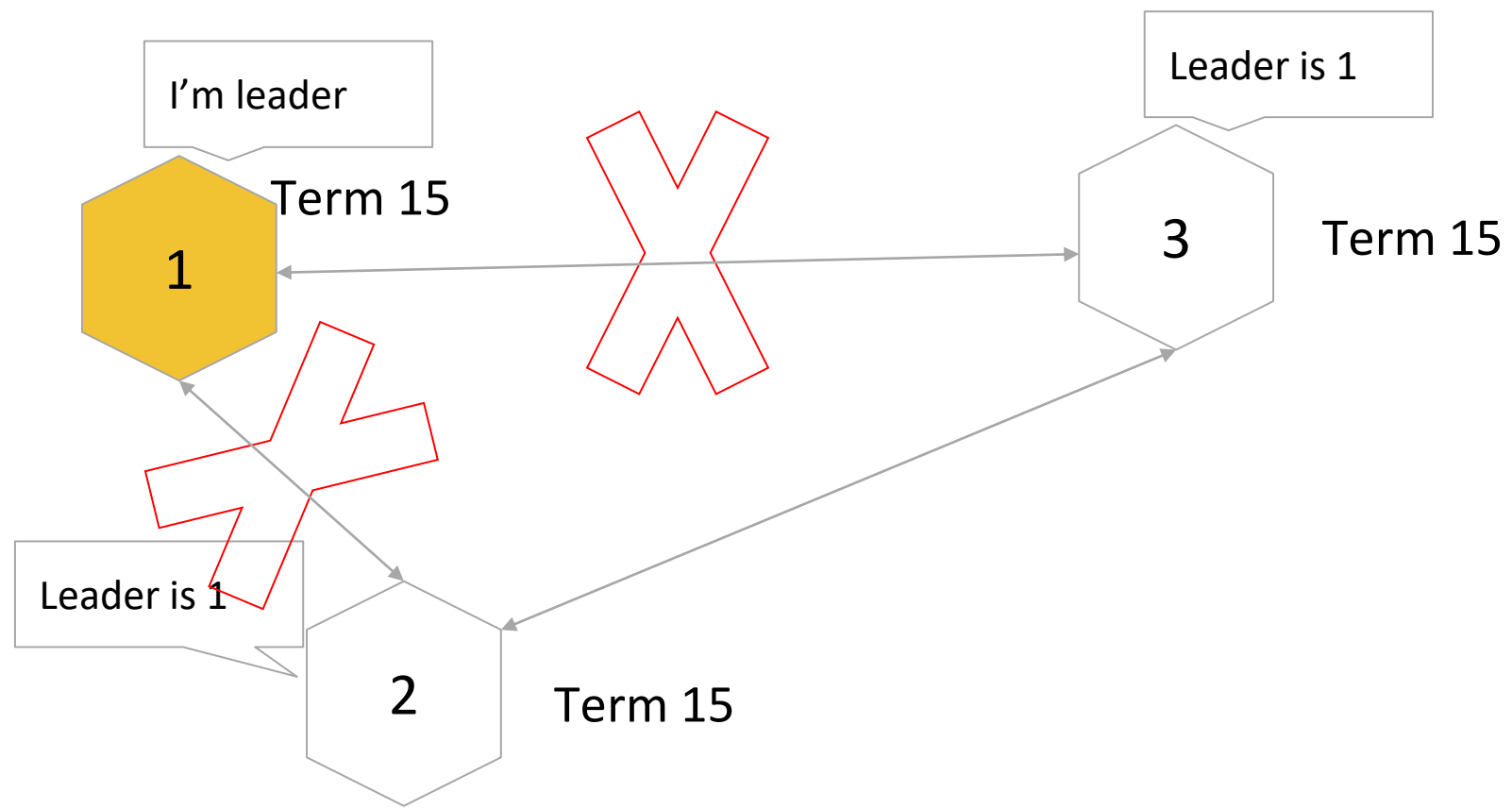
- Для записи в кластер необходимо обратиться к текущему лидеру
- Возможность существования нескольких нод, считающих себя лидерами одновременно (в разных term)
- Не более одного из “лидеров” сможет на самом деле произвести синхронную запись



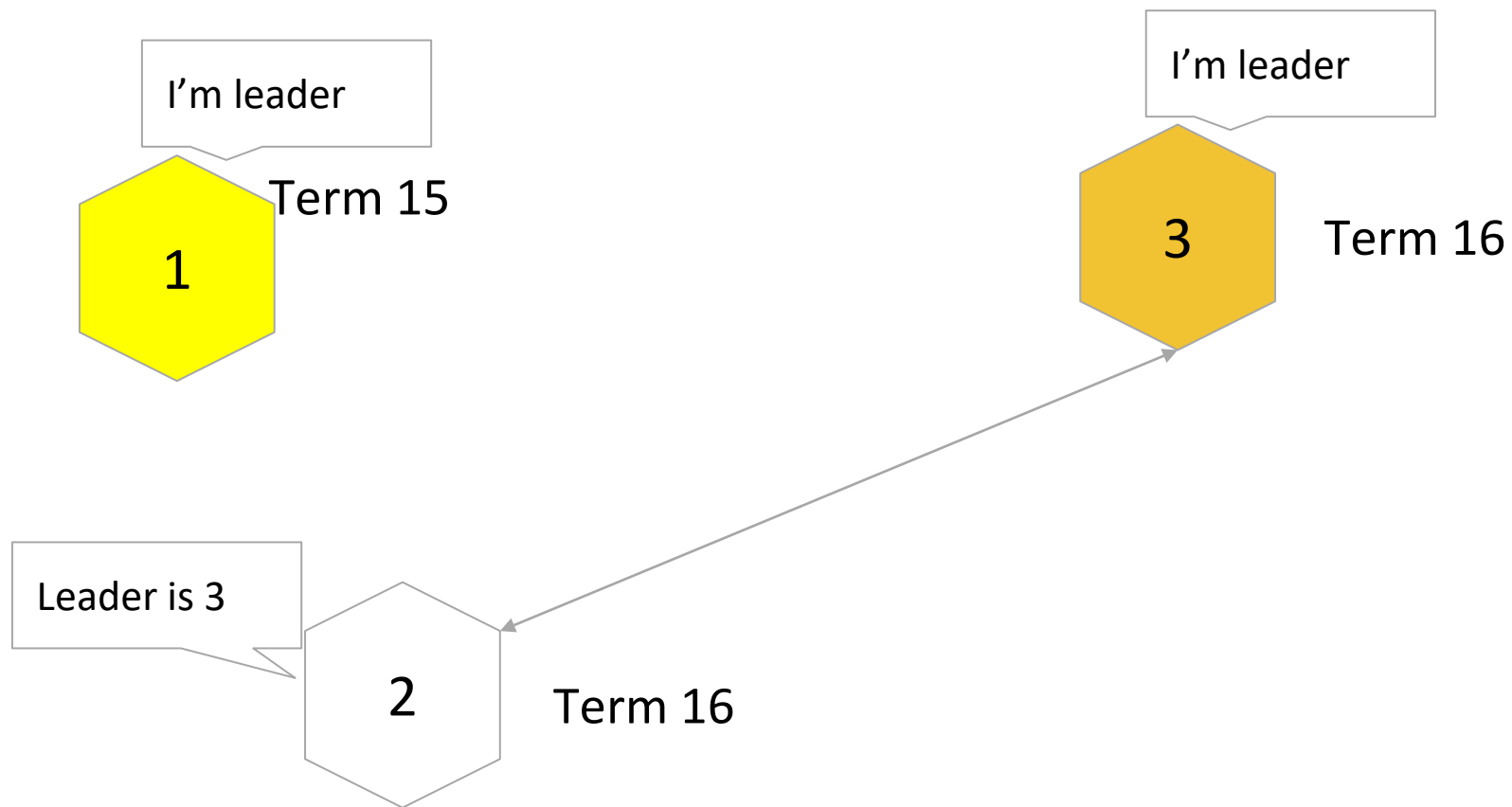
# Fencing: Пример №1



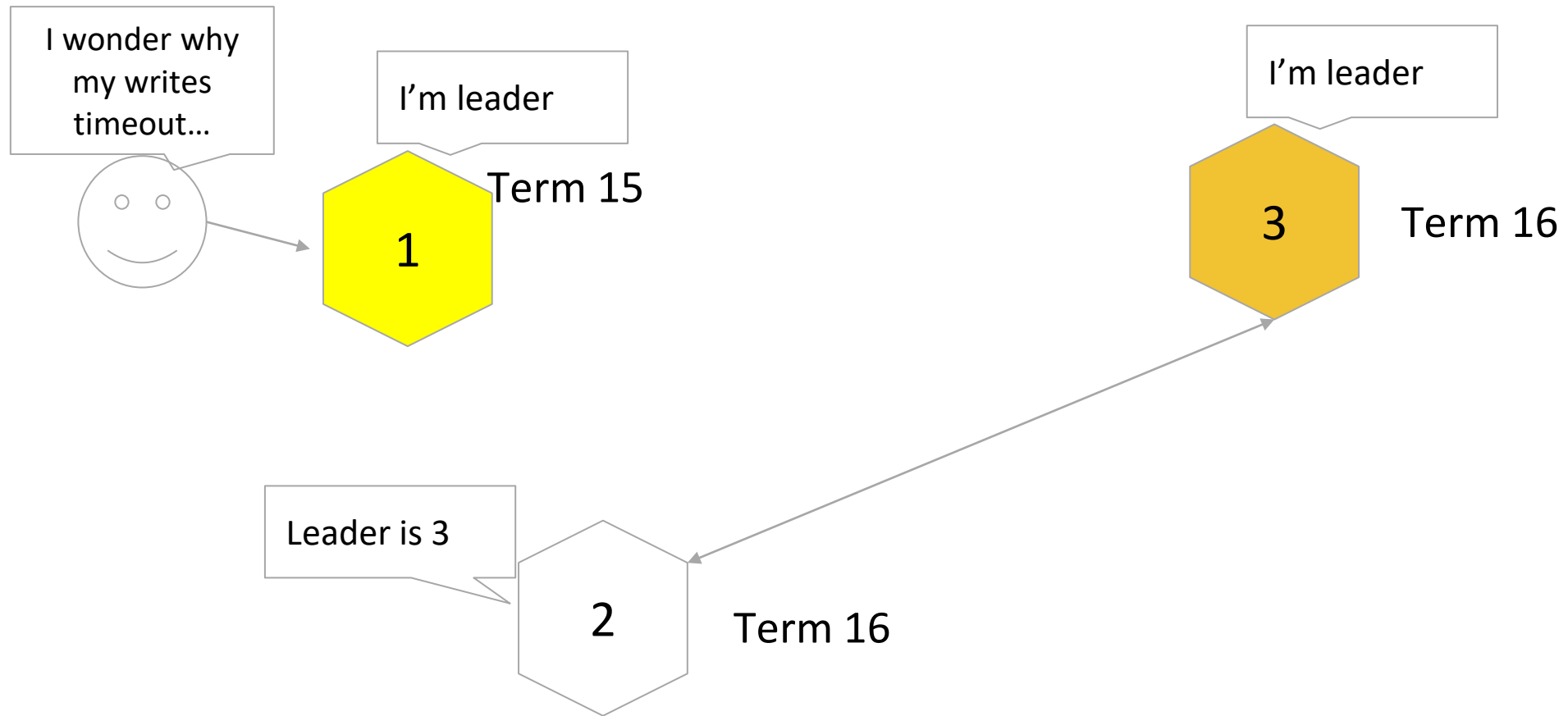
# Fencing: Пример №1



# Fencing: Пример №1



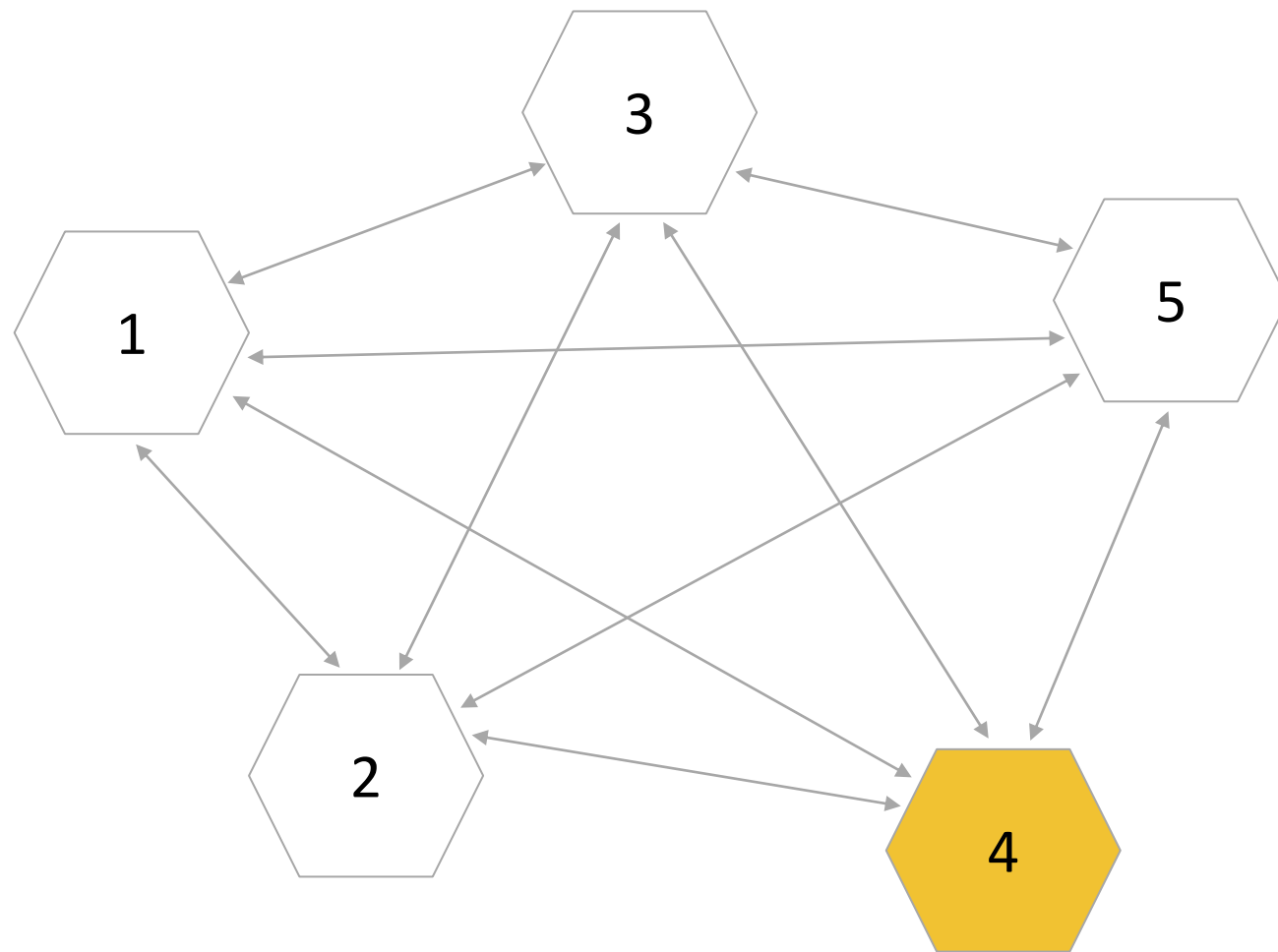
# Fencing: Пример №1



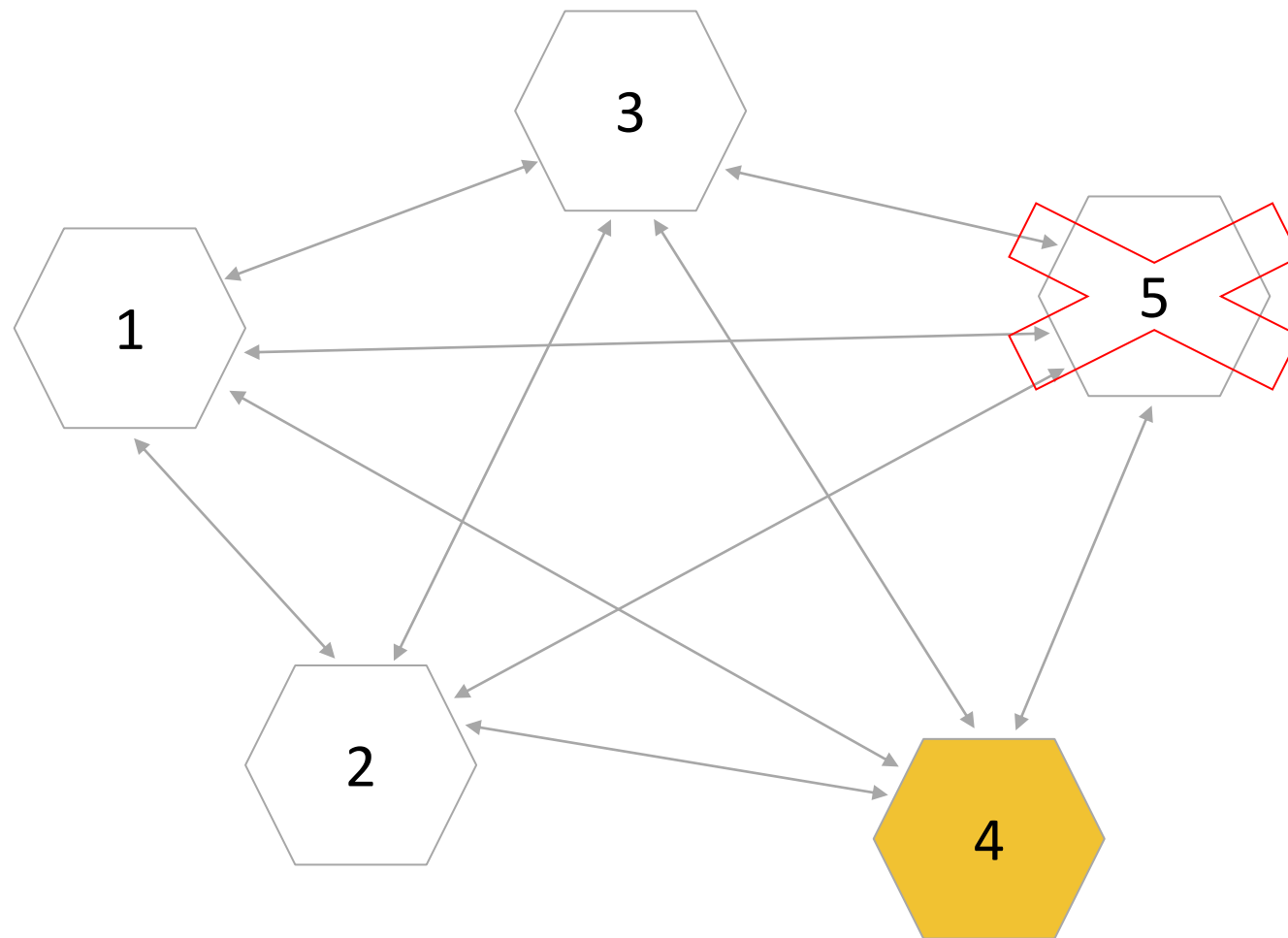
# Fencing: Проблема №2

- При наличии pre-vote можно получить “заблокированный” кластер

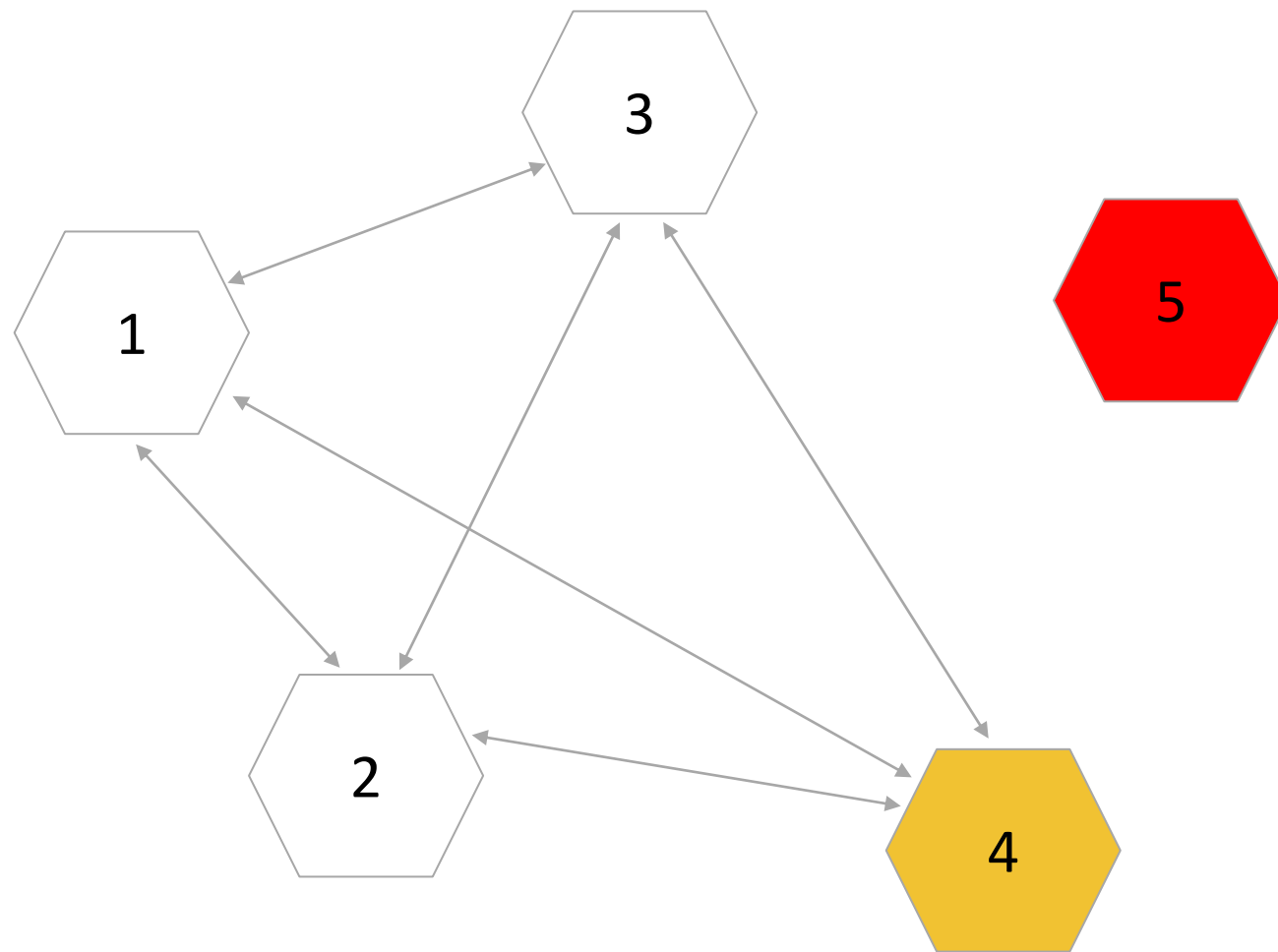
# Fencing: Пример №2



# Fencing: Пример №2

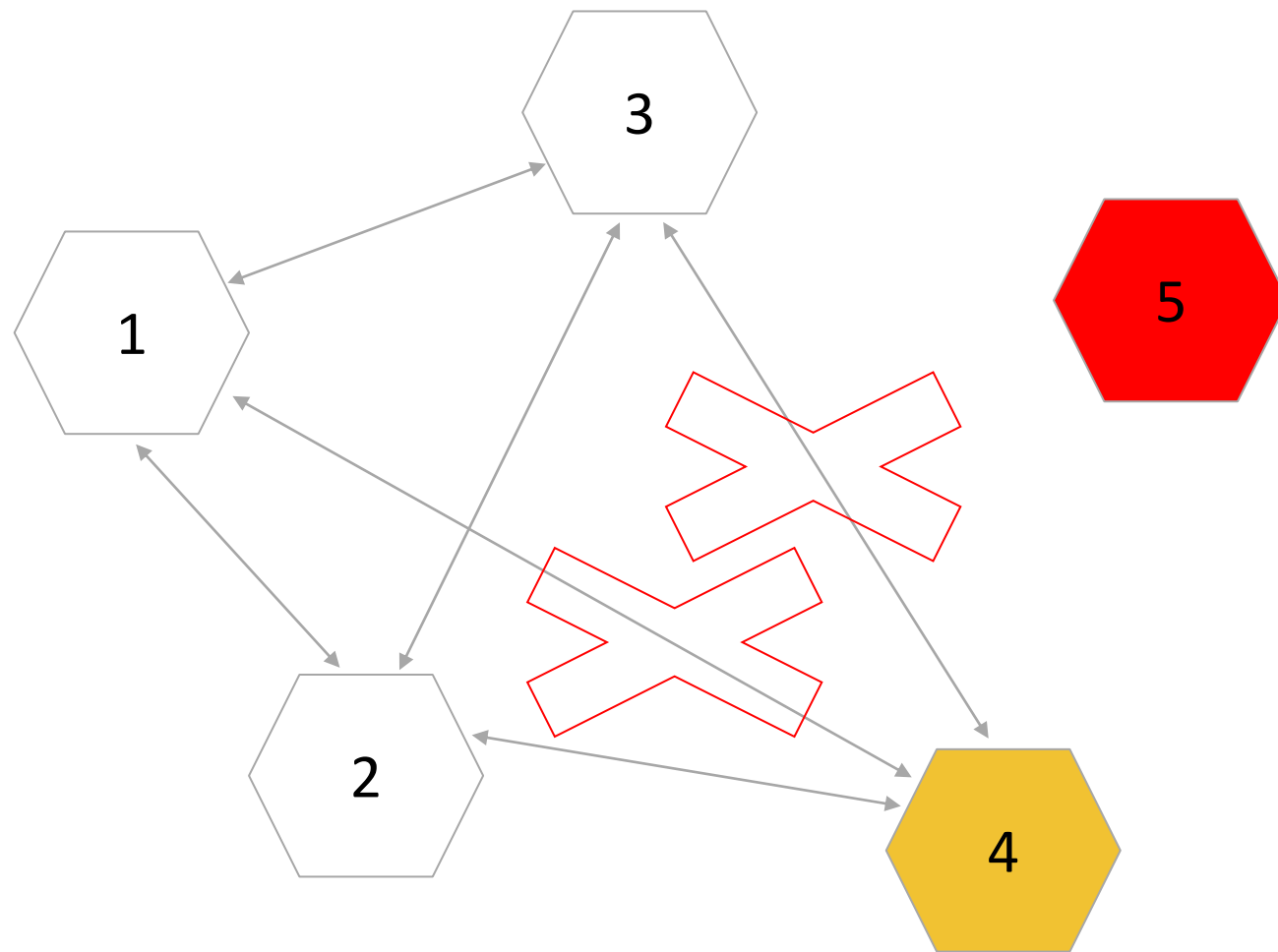


# Fencing: Пример №2

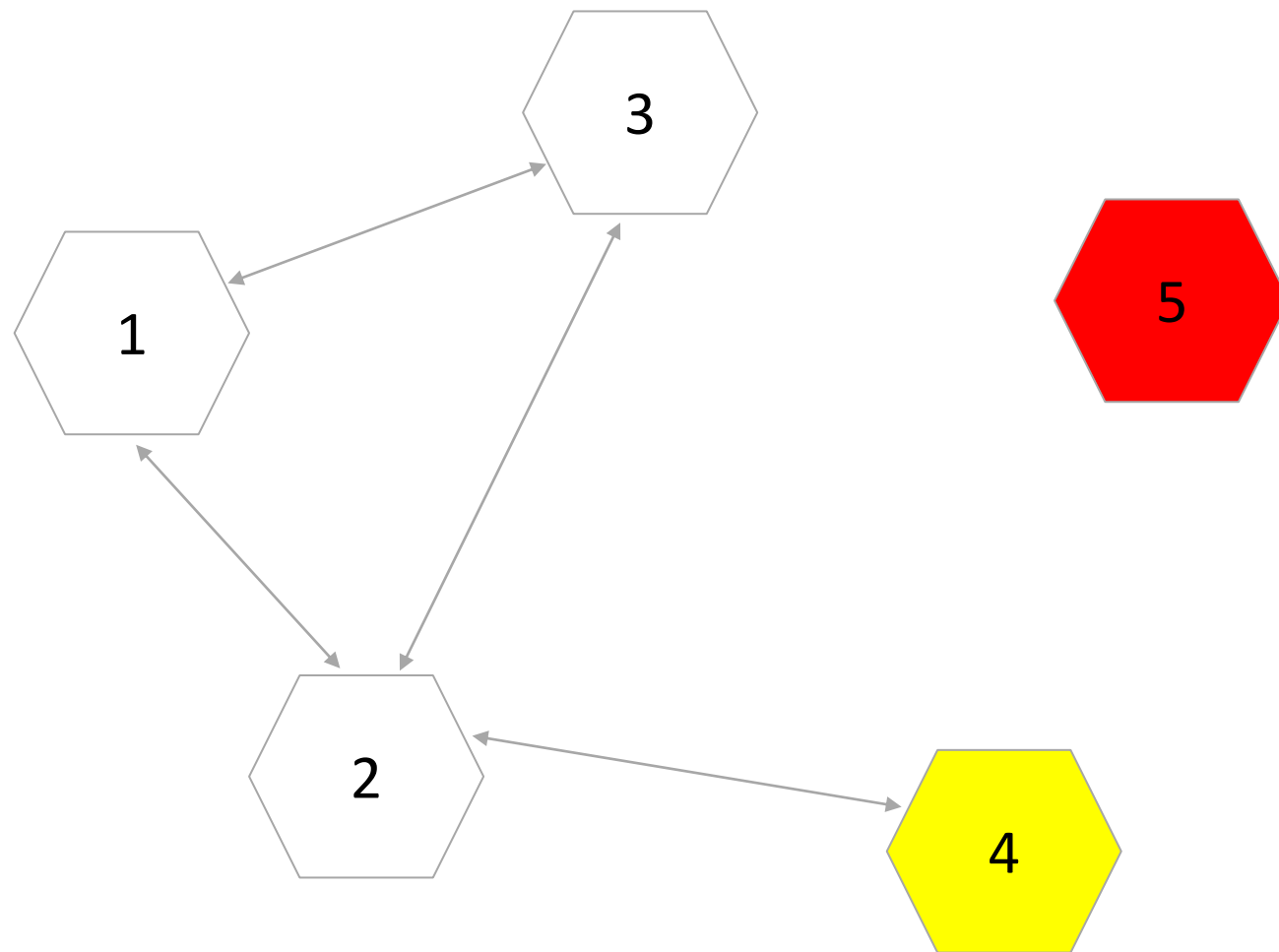




# Fencing: Пример №2



# Fencing: Пример №2



# Fencing: Способы решения

- Для определения текущего лидера запрашивать у всех нод кластера текущий терм и состояние ноды (лидер или нет)
- Переложить определение лидера на внешнюю систему, гарантирующую его единственность, и довериться ей
- Доработать RAFT в Tarantool

# Fencing: Способы решения

1. Для определения текущего лидера запрашивать у всех нод кластера текущий терм и состояние ноды (лидер или нет)
2. Переложить определение лидера на внешнюю систему, гарантирующую его единственность, и довериться ей
3. Доработать RAFT в Tarantool

# Fencing: Способы решения

1. Для определения текущего лидера запрашивать у всех нод кластера текущий терм и состояние ноды (лидер или нет)
  - Не решает проблему заблокированного кластера
2. Переложить определение лидера на внешнюю систему, гарантирующую его единственность, и довериться ей
3. Доработать RAFT в Tarantool

# Fencing: Способы решения

1. Для определения текущего лидера запрашивать у всех нод кластера текущий терм и состояние ноды (лидер или нет)
  - Не решает проблему заблокированного кластера
  - Увеличивается время недоступности на запись
2. Переложить определение лидера на внешнюю систему, гарантирующую его единственность, и довериться ей
3. Доработать RAFT в Tarantool

# Fencing: Способы решения

1. Для определения текущего лидера запрашивать у всех нод кластера текущий терм и состояние ноды (лидер или нет)
  - Не решает проблему заблокированного кластера
  - Увеличивается время недоступности на запись
  - Требуется реализация данного решения во всех клиентах
2. Переложить определение лидера на внешнюю систему, гарантирующую его единственность, и довериться ей
3. Доработать RAFT в Tarantool

# Fencing: Способы решения

1. Для определения текущего лидера запрашивать у всех нод кластера текущий терм и состояние ноды (лидер или нет)
2. **Переложить определение лидера на внешнюю систему, гарантирующую его единственность, и довериться ей**
3. Доработать RAFT в Tarantool



# Fencing: Способы решения

1. Для определения текущего лидера запрашивать у всех нод кластера текущий терм и состояние ноды (лидер или нет)
2. **Переложить определение лидера на внешнюю систему, гарантирующую его единственность, и довериться ей**
  - Доверяем решение проблемы третьей системе
3. Доработать RAFT в Tarantool

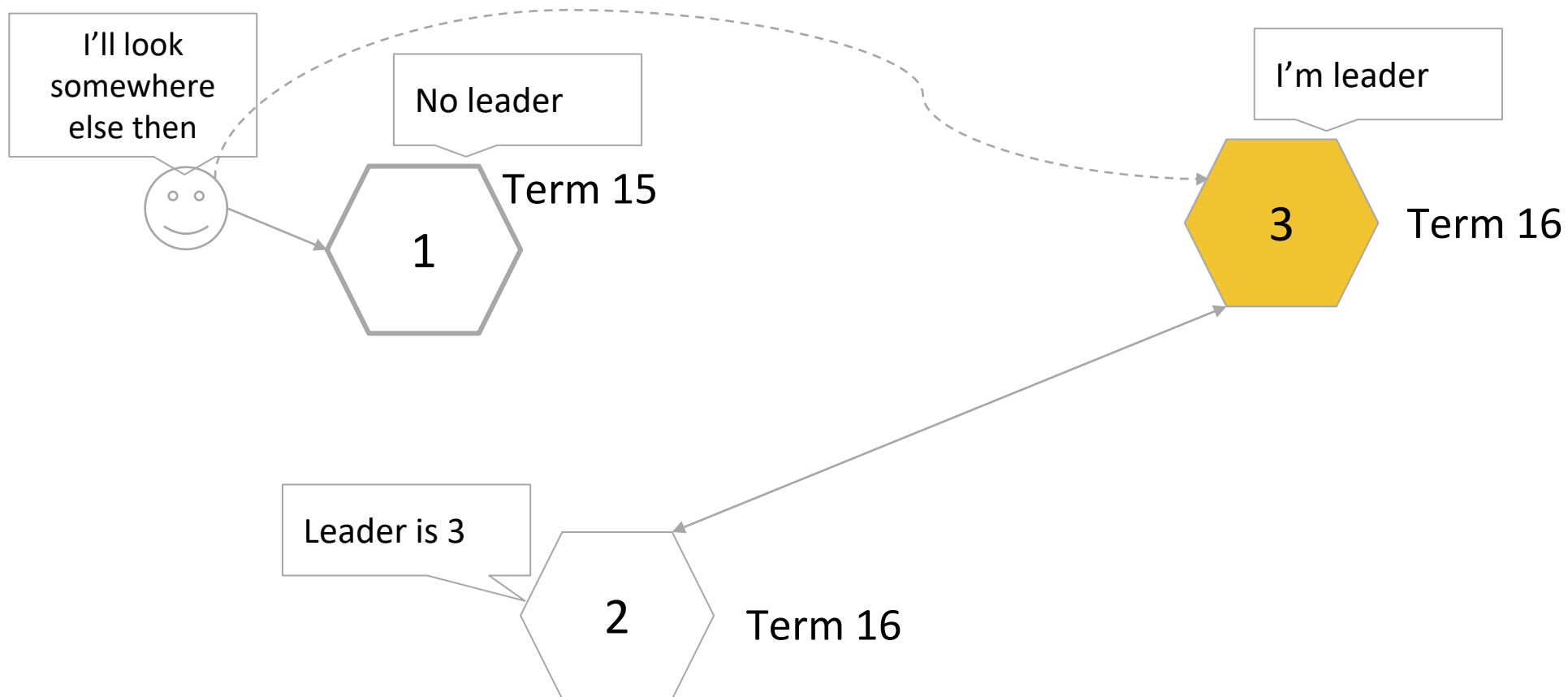
# Fencing: Способы решения

1. Для определения текущего лидера запрашивать у всех нод кластера текущий терм и состояние ноды (лидер или нет)
2. **Переложить определение лидера на внешнюю систему, гарантирующую его единственность, и довериться ей**
  - Доверяем решение проблемы третьей системе
  - Еще одна точка отказа
3. Доработать RAFT в Tarantool

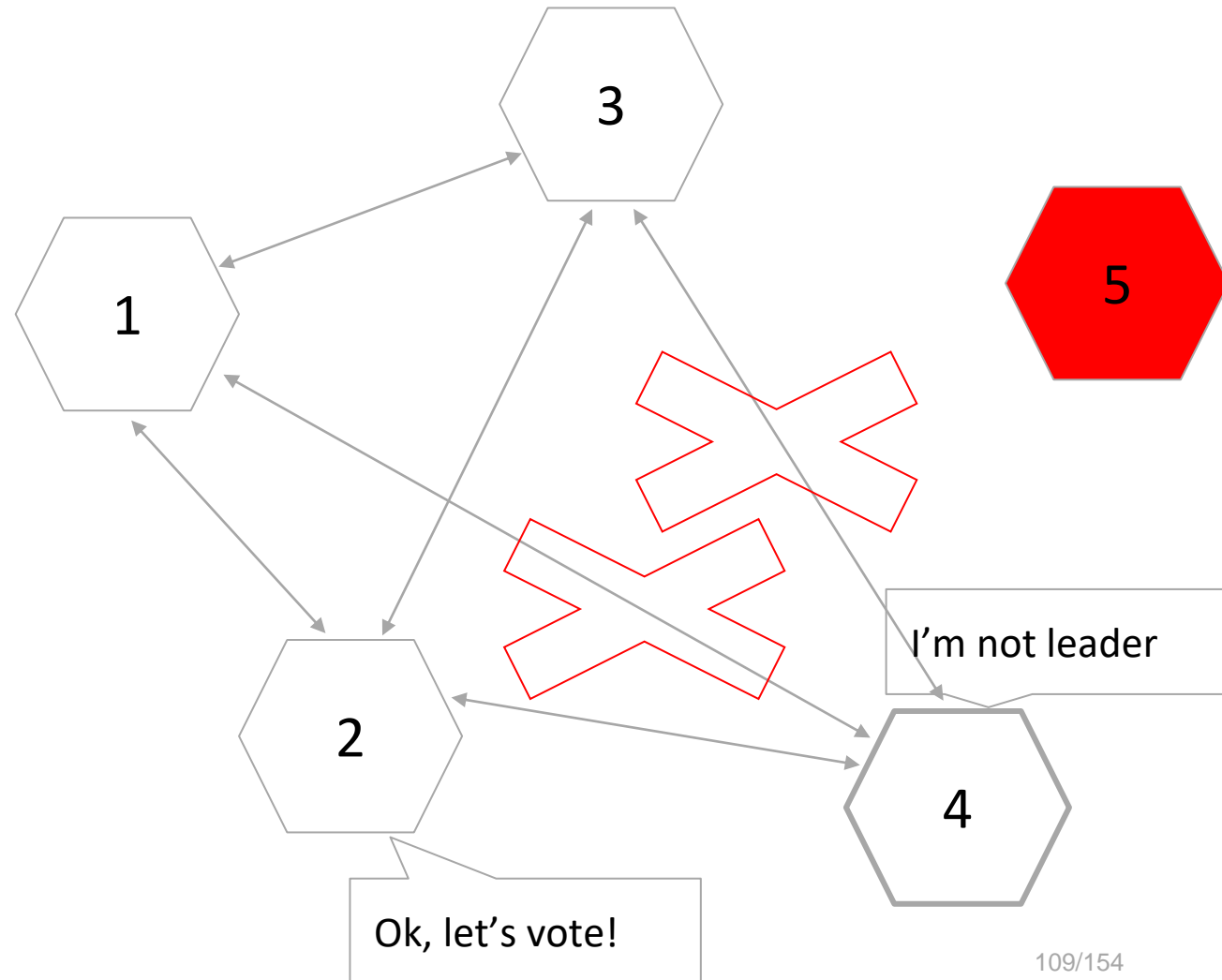
# Fencing: Способы решения

1. Для определения текущего лидера запрашивать у всех нод кластера текущий терм и состояние ноды (лидер или нет)
  - Переложить определение лидера на внешнюю систему, гарантирующую его единственность, и довериться ей
  - **Доработать RAFT в Tarantool**

# Fencing: Как хочется



# Как хочется

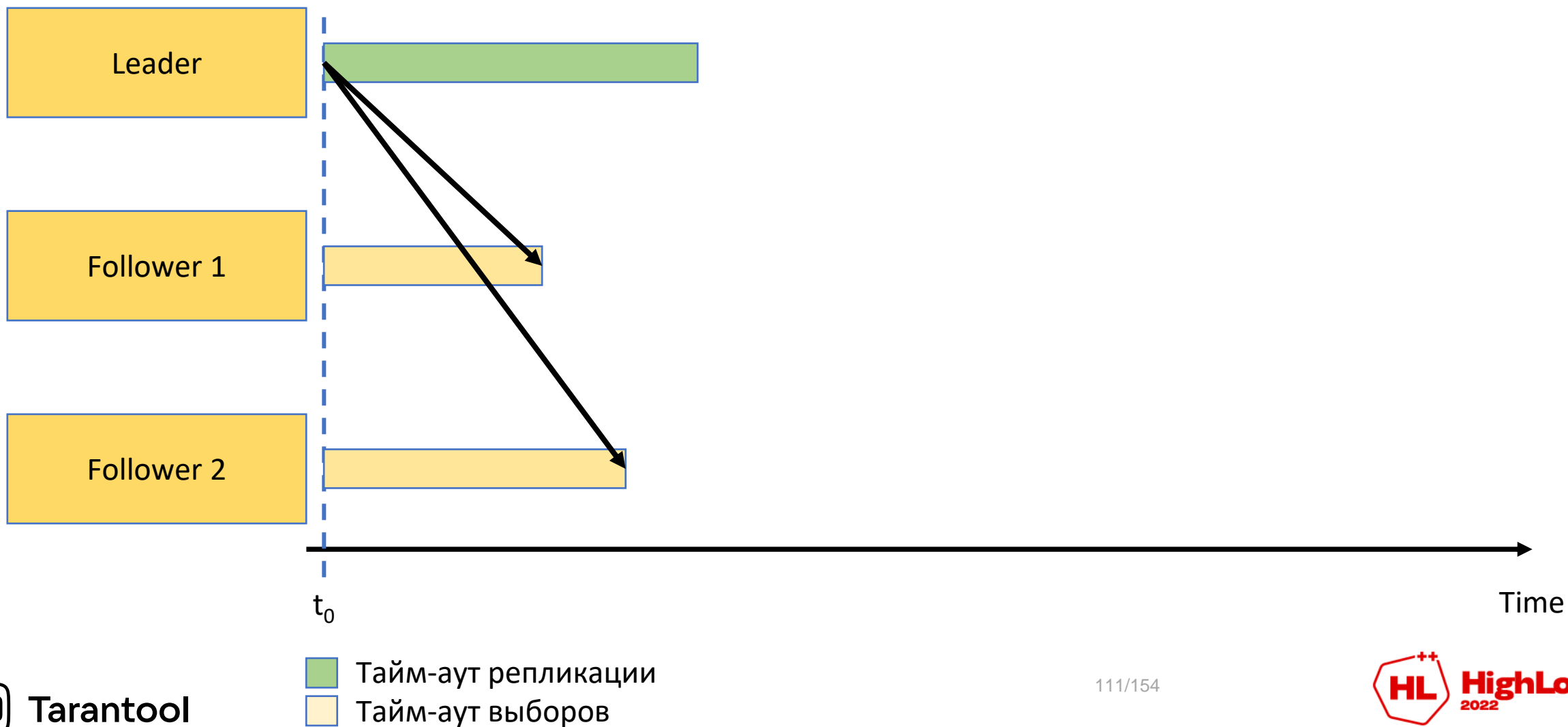


# Fencing: Решение

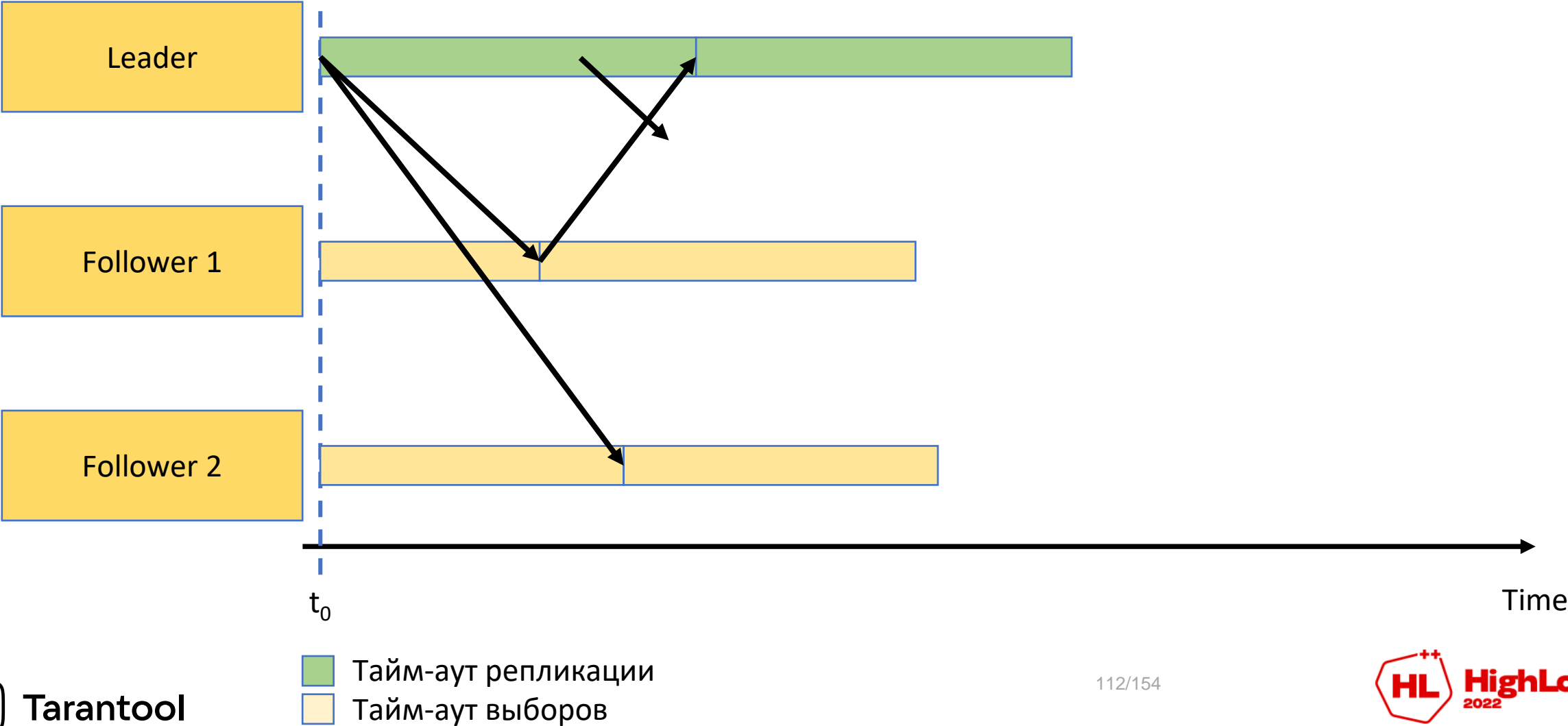
Нода, являющаяся лидером, следит за обрывом соединений к/от реплик

Если количество живых соединений стало меньше, чем текущий кворум, лидер снимает с себя полномочия

# Fencing: Решение

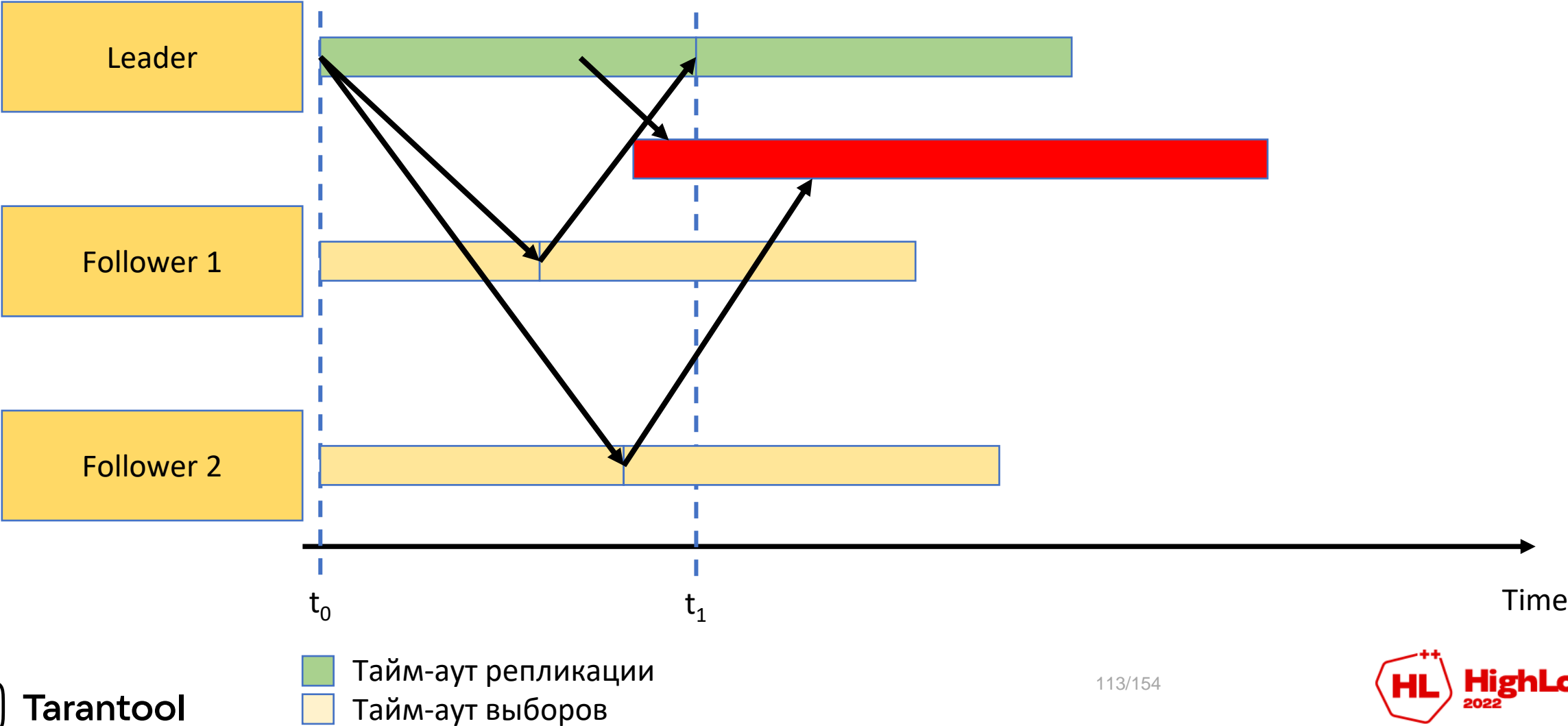


# Fencing: Решение

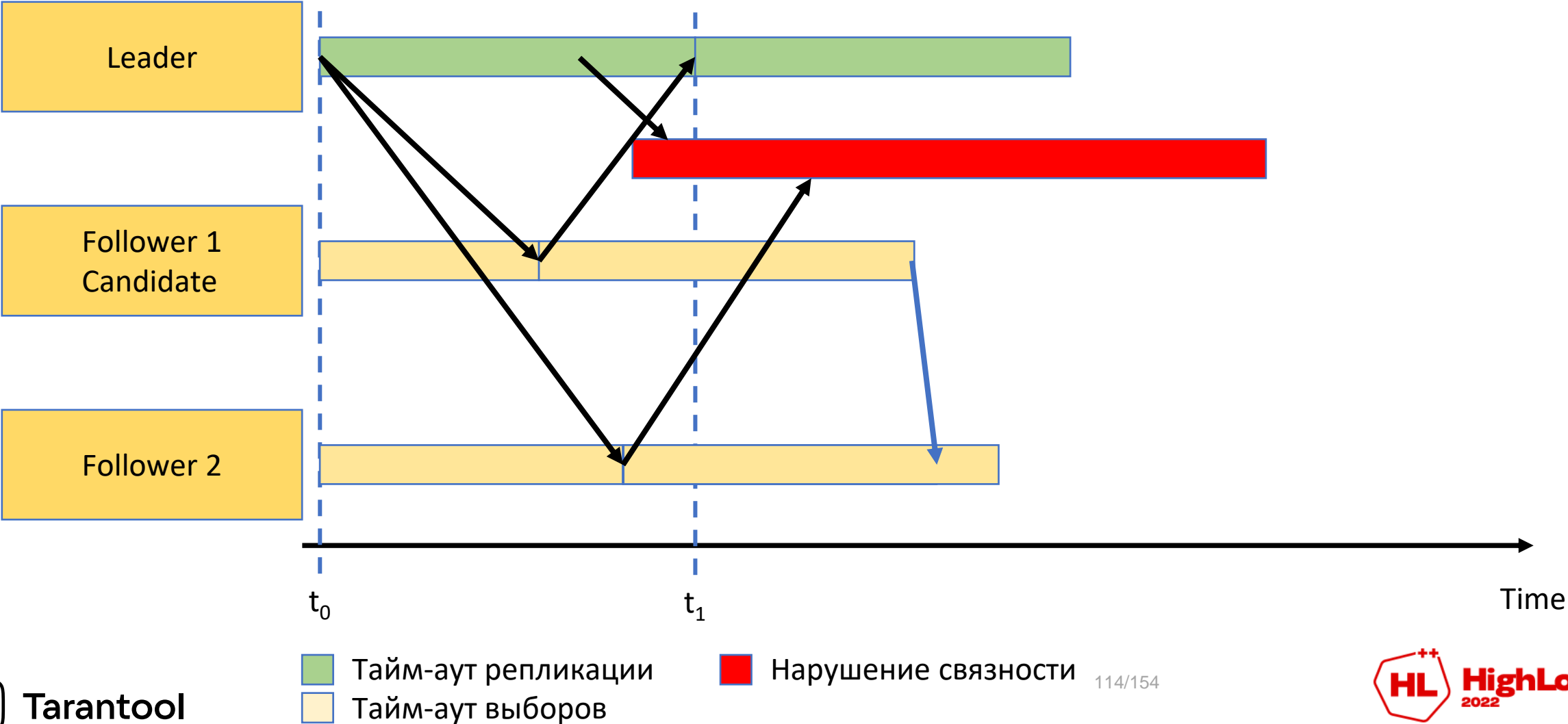




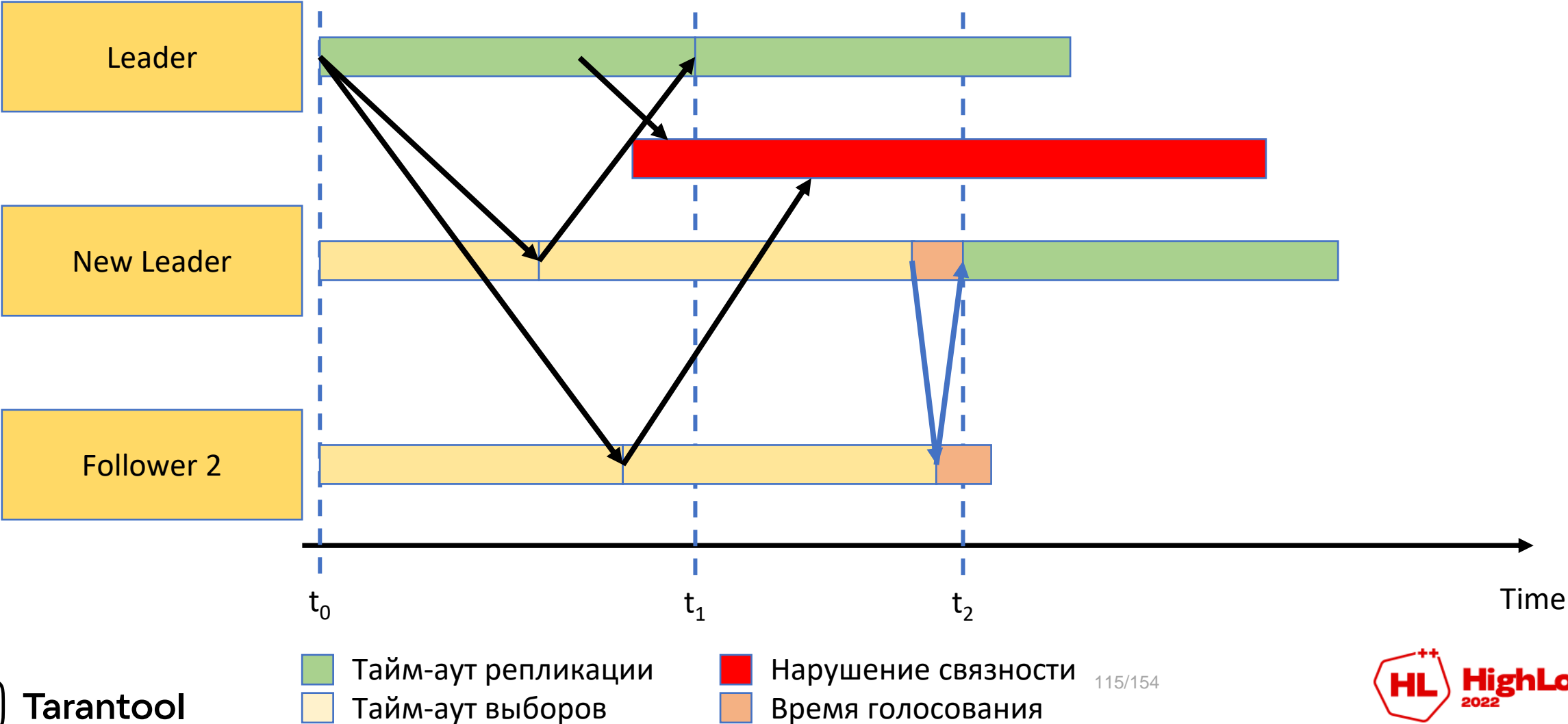
# Fencing: Решение



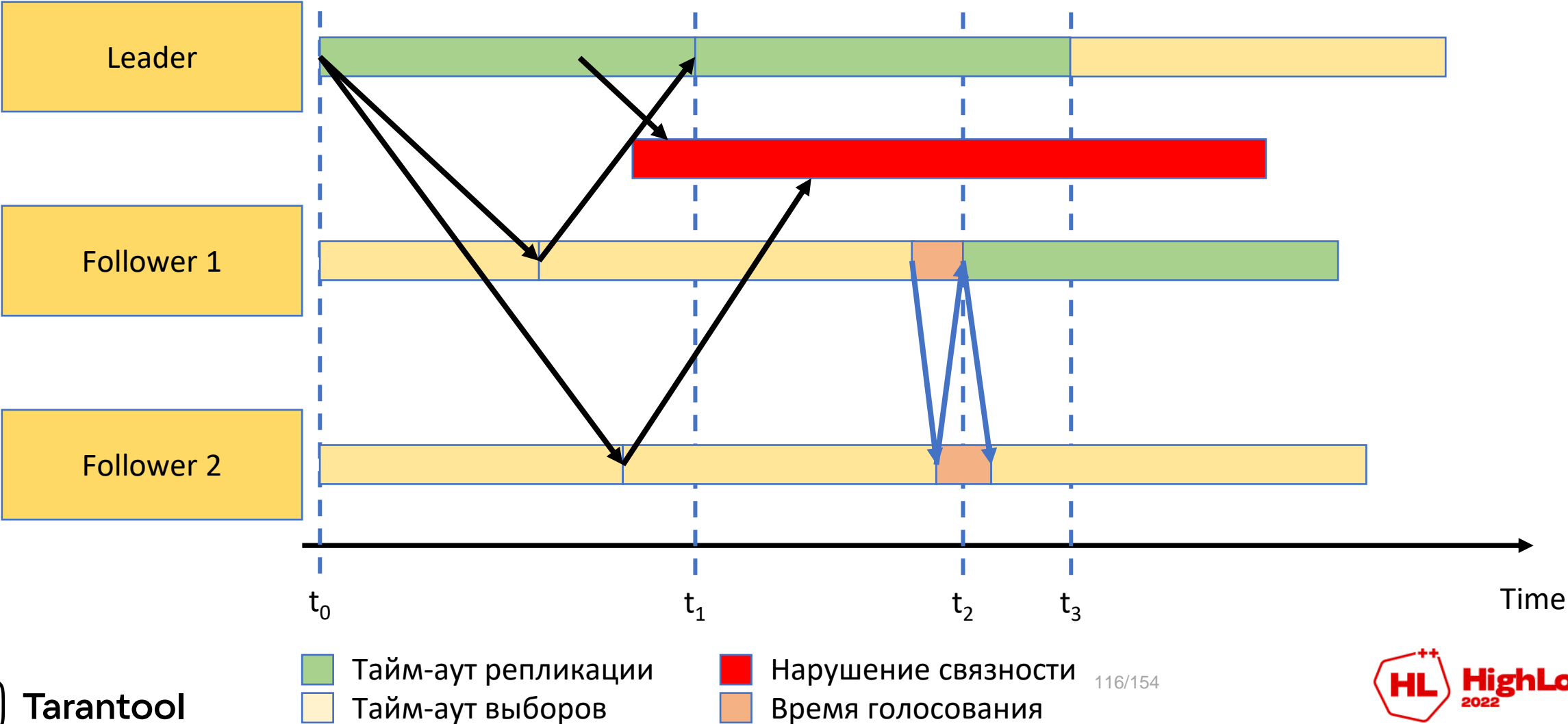
# Fencing: Решение



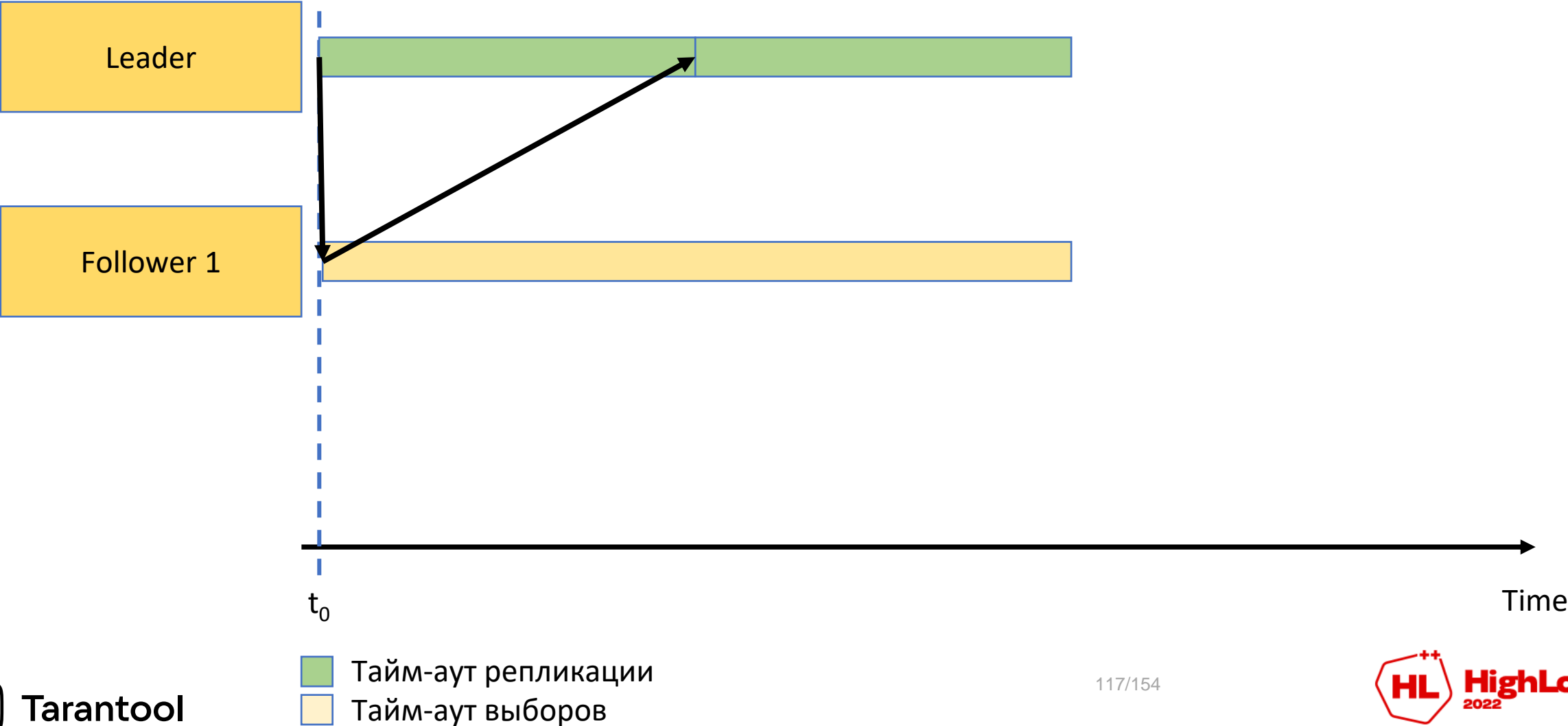
# Fencing: Решение



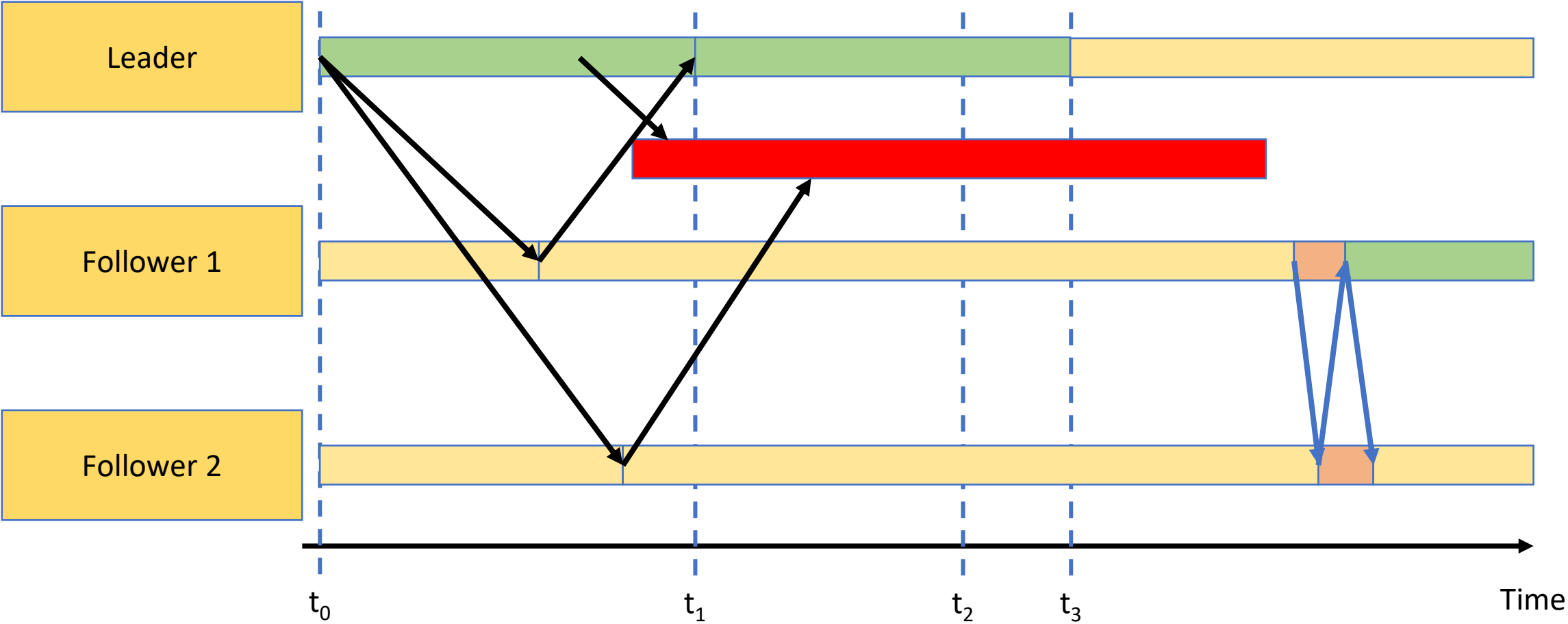
# Fencing: Решение



# Fencing: Решение



# Fencing: Решение



- Тайм-аут репликации

Тайм-аут выборов

Нарушение связности

Время голосования
- 118/154

# Fencing: Решение

Три варианта применения Fencing:

1. Fencing выключен
2. Слабый fencing – более быстрый фейловер, но можно увидеть двух лидеров
3. Сильный fencing – точно не будет двух лидеров

# Ожидали от Raft

- Один лидер на терм == один лидер в данный момент времени
- Лидер постоянен
- Быстрая смена лидера



# Получили

- Один лидер на терм == один лидер в данный момент времени – Fencing
- Лидер постоянен – Pre-Vote, Fencing
- Смены лидера – Split-Vote detection позволил ускорить выборы ~10 раз

Обратная связь  
и комментарии по докладу  
по ссылке



Tarantool